last time

preview: OS enforces restrictions

user/group IDs

access control list idea list of who can access what

chmod rwx for owner/one group/everyone else

more general access control lists

superuser/root = always has permissionbut still goes through OS to do things

set-user-ID: controlled access to special functionality mark specific programs to have extra access

reminder: warmup

problem: existing executable breaks when .so file updated

if .a instead of .so:

executable includes machine code from *old version* of library because linking process includes .a file

bar.h has both prototypes

would not compile+link in C (but not C++) — can't have two prototypes for one function (in C++ — would need to *define* both versions of function)

all: executable

```
foo.c (should be foo.o): foo.h
  gcc -Wall -c foo.c
```

main.c (should be main.o): foo.h
gcc -Wall -c main.c

executable: foo.c main.c (should be foo.o main.o) gcc -Wall -o executable foo.c main.c (should be foo.o main.o)

clean:

rm --force main.o foo.o executable



if foo.h more recent than main.c, for builds main.o from main.c (compiles main.c once) builds foo.o from foo.c (A) executable need not be updated (if main.c/foo.c old) (C)

if foo.c or main.c more recent than builds executable from foo.c, main (compile main.c once)



if foo.h more recent than main.c, for builds main.o from main.c (compiles main.c once) builds foo.o from foo.c (A) executable need not be updated (if main.c/foo.c old) (C)

if foo.c or main.c more recent than builds executable from foo.c, main (compile main.c once)



if foo.h more recent than main.c, for builds main.o from main.c (compiles main.c once) builds foo.o from foo.c (A) executable need not be updated (if main.c/foo.c old) (C)

if foo.c or main.c more recent than o builds executable from foo.c, main (compile main.c once)



if foo.h more recent than main.c, for builds main.o from main.c (compiles main.c once) builds foo.o from foo.c (A) executable need not be updated (if main.c/foo.c old) (C)

if foo.c or main.c more recent than builds executable from foo.c, main (compile main.c once)

Q3

logo.h: logo.bmp

xxd reads logo.bmp — is prerequisite
xxd writes logo.h — is target

main.o: logo.h graphics.h

don't need command because it's in %.o: %.c rule (but not error)

don't need main.c because it's in %.o: %.c rule (but not error)

do need logo.h, graphics.h

shouldn't have stdio.h since it's not in this directory
 could have /usr/include/stdio.h or similar

Q6

want two files: one r/w for A, B only one r/w for A, C only

group with A: owned by B, user rw, group-with-A rw, other no access owner by C, user rw, group-with-A rw, other no access

group with B + group with C: owned by A, user rw, group-with-B rw, other no access owner by A, user rw, group-with-C rw, other no access things programs on portal shouldn't do

read other user's files

modify OS's memory

read other user's data in memory

hang the entire system

things programs on portal shouldn't do

read other user's files

modify OS's memory

read other user's data in memory

hang the entire system

privileged operation: problem

how can hardware (HW) plus operating system (OS) allow: read your own files from hard drive

but disallow:

read others files from hard drive

some ideas

OS tells HW 'okay' parts of hard drive before running program code

complex for hardware and for OS

some ideas

OS tells HW 'okay' parts of hard drive before running program code

complex for hardware and for OS

OS verifies your program's code can't do bad hard drive access no work for HW, but complex for OS may require compiling differently to allow analysis

some ideas

OS tells HW 'okay' parts of hard drive before running program code

complex for hardware and for OS

OS verifies your program's code can't do bad hard drive access no work for HW, but complex for OS may require compiling differently to allow analysis

OS tells HW to only allow OS-written code to access hard drive that code can enforce only 'good' accesses requires program code to call OS routines to access hard drive relatively simple for hardware

kernel mode

extra one-bit register: "are we in *kernel mode*" other names: privileged mode, supervisor mode, ...

```
not in kernel mode = user mode
```

```
certain operations only allowed in kernel mode 
privileged instructions
```

```
example: talking to any I/O device
```

what runs in kernel mode?

system boots in kernel mode

OS switches to user mode to run program code

next topic: when does system switch back to kernel mode? how does OS tell HW where the (trusted) OS code is?

hardware + system call interface

applications $+$ libraries					
user-mode hardware interface (limited)	system call interface				
	kernel part of OS that runs in kernel mode				
	kernel-mode hardware interface (complete)				
hardware					



controlled entry to kernel mode (1)

special instruction: "make system call"
similar idea as call instruction — jump to function elsewhere
(and allow that function to return later)

runs OS code in kernel mode at location specified earlier

OS sets up at boot

location can't be changed without privilieged instrution

controlled entry to kernel mode (2)

OS needs to make specified location:

figure out what operation the program wants calling convention, similar to function arguments + return value

be "safe" — not allow the program to do 'bad' things example: checks whether current program is allowed to read file before reading it requires exceptional care — program can try weird things





system call terminology

some inconsistency:

system call = event of entering kernel mode on request?

system call = whole porcess from beginning to end?

same issue as with 'function call' is it just starting the function, or the whole time the function runs?

keeping permissions?

which of the following would still be secure?

A. performing authorization checks in the standard library in addition to system call handlers

B. performing authorization checks in the standard library instead of system call handlers

C. making the user ID a system call argument rather than storing it persistently in the OS's memory

Linux x86-64 system calls

special instruction: syscall

runs OS specified code in kernel mode

Linux syscall calling convention

before syscall:

- %rax system call number
- %rdi, %rsi, %rdx, %r10, %r8, %r9 args

after syscall:

%rax — return value

on error: %rax contains -1 times "error number"

almost the same as normal function calls

Linux x86-64 hello world

syscall

```
.globl start
.data
hello_str: .asciz "Hello, World!\n"
.text
start:
  movg $1, %rax # 1 = "write"
  movg $1, %rdi # file descriptor 1 = stdout
  mova $hello str. %rsi
  movg $15, %rdx # 15 = strlen("Hello, World!\n")
  svscall
  movg $60, %rax # 60 = exit
  movq $0, %rdi
```

approx. system call handler

```
sys_call_table:
    .quad handle_read_syscall
    .quad handle_write_syscall
    // ...
```

```
handle syscall:
    ... // save old PC, etc.
    pushq %rcx // save registers
    pusha %rdi
    . . .
    call *sys call table(,%rax,8)
    . . .
    popq %rdi
    popq %rcx
    return from exception
```

Linux system call examples

mmap, brk — allocate memory

fork — create new process

execve — run a program in the current process

_exit — terminate a process

open, read, write — access files

socket, accept, getpeername — socket-related

Linux system call examples

mmap, brk — allocate memory

fork — create new process

execve — run a program in the current process

_exit — terminate a process

open, read, write — access files

socket, accept, getpeername — socket-related









system call wrappers

library functions to not write assembly:

```
open:
    movg $2, %rax // 2 = sys_open
    // 2 arguments happen to use same registers
    svscall
    // return value in %eax
    cmp $0, %rax
    il has error
    ret
has error:
    neg %rax
    movq %rax, errno
    movq $-1, %rax
    ret
```

system call wrappers

library functions to not write assembly:

```
open:
    movg $2, %rax // 2 = sys_open
    // 2 arguments happen to use same registers
    svscall
    // return value in %eax
    cmp $0, %rax
    il has error
    ret
has error:
    neg %rax
    movq %rax, errno
    movg $-1, %rax
    ret
```

system call wrapper: usage

```
/* unistd.h contains definitions of:
    O_RDONLY (integer constant), open() */
#include <unistd.h>
int main(void) {
  int file descriptor:
  file descriptor = open("input.txt", O RDONLY);
  if (file_descriptor < 0) {</pre>
      printf("error: %s\n", strerror(errno));
      exit(1):
  }
  result = read(file descriptor, ...);
  . . .
```

system call wrapper: usage

```
/* unistd.h contains definitions of:
    O_RDONLY (integer constant), open() */
#include <unistd.h>
int main(void) {
  int file descriptor:
  file_descriptor = open("input.txt", 0_RDONLY);
  if (file_descriptor < 0) {</pre>
      printf("error: %s\n", strerror(errno));
      exit(1):
  }
  result = read(file descriptor, ...);
  . . .
```

strace hello_world (1)

strace — Linux tool to trace system calls

strace hello_world (2)

#include <stdio.h>
int main() { puts("Hello, World!"); }

```
when statically linked:
execve("./hello_world", ["./hello_world"], 0x7ffeb4127f70 /* 28 vars */)
                                         =
                                           0
brk(NULL)
                                         = 0 \times 22 f 8000
brk(0x22f91c0)
                                         = 0x22f91c0
arch_prctl(ARCH_SET_FS, 0x22f8880)
                                         = 0
uname({sysname="Linux", nodename="reiss-t3620", ...}) = 0
readlink("/proc/self/exe", "/u/cr4bd/spring2023/cs3130/slide"..., 4096)
                                         = 57
brk(0x231a1c0)
                                         = 0x231a1c0
brk(0x231b000)
                                         = 0x231b000
access("/etc/ld.so.nohwcap", F_OK)
                                         = -1 ENOENT (No such file or
                                                       directorv)
fstat(1, {st mode=S IFCHR|0620, st rdev=makedev(136, 4), ...}) = 0
write(1, "Hello, World!\n", 14)
                                         = 14
exit_group(0)
                                           ?
                                         =
                                                                          32
     . . . . . .
```

aside: what are those syscalls?

execve: run program

brk: allocate heap space

- arch_prctl(ARCH_SET_FS, ...): thread local storage pointer may make more sense when we cover concurrency/parallelism later
- uname: get system information
- readlink of /proc/self/exe: get name of this program
- access: can we access this file [in this case, a config file]?
- fstat: get information about open file
- exit_group: variant of exit

strace hello_world (2)

```
#include <stdio.h>
int main() { puts("Hello, World!"); }
when dynamically linked:
execve("./hello_world", ["./hello_world"], 0x7ffcfe91d540 /* 28 vars */)
                                        =
brk(NULL)
                                        = 0x55d6c351b000
. . .
openat(AT_FDCWD, "/etc/ld.so.cache", O_RDONLY|O_CLOEXEC) = 3
fstat(3, {st_mode=S_IFREG|0644, st_size=196684, ...}) = 0
mmap(NULL, 196684, PROT_READ, MAP_PRIVATE, 3, 0) = 0x7f7a62dd3000
close(3)
                                        = 0
access("/etc/ld.so.nohwcap", F OK) = -1 ENOENT (No such file or director)
openat(AT_FDCWD, "/lib/x86_64-linux-gnu/libc.so.6", 0 RDONLY|0 CLOEXEC) = 3
read(3, "177ELF(2)(1)(3)(0)(0)(0)(0)(0)(3)(0)(1)(0)(0)(0), 832) = 832
. . .
```

```
close(3) = 0
write(1, "Hello, World!\n", 14) = 14
exit_group(0)
+++ exited with 0 +++
```

backup slides

crash timeline timeline





keyboard input timeline



```
handle_timer_interrupt:
   save_old_pc save_pc
   movq %r15, save_r15
   /* key press here */
```

```
movq %r14, save_r14
```

. . .





interrupt disabling

CPU supports disabling (most) interrupts

interrupts will wait until it is reenabled

CPU has extra state:

are interrupts enabled? is keyboard interrupt pending? is timer interrupt pending?

```
handle timer interrupt:
 /* interrupts automatically disabled here */
 movq %rsp, save_rsp
  save old pc save pc
 /* kev press here */
  impIfFromKernelMode skip exception stack
 movg current exception stack, %rsp
skip set kernel stack:
  pushq save rsp
  pushq save_pc
  enable intterupts2
  pusha %r15
  . . .
```

```
/* interrupt happens here! */
```

```
. . .
```

```
handle timer interrupt:
 /* interrupts automatically disabled here */
 movq %rsp, save_rsp
  save old pc save pc
 /* kev press here */
  impIfFromKernelMode skip exception stack
 movg current exception stack, %rsp
skip_set_kernel_stack:
  pushq save rsp
  pushq save pc
  enable intterupts2
  pusha %r15
```

```
• • •
```

```
/* interrupt happens here! */
```

```
handle timer interrupt:
 /* interrupts automatically disabled here */
 movq %rsp, save_rsp
  save old pc save pc
 /* kev press here */
  impIfFromKernelMode skip_exception_stack
 movg current_exception_stack, %rsp
skip set kernel stack:
  pushq save rsp
  pushq save_pc
  enable intterupts2
  pusha %r15
  . . .
```

/* interrupt happens here! */

disabling interrupts

automatically disabled when exception handler starts

also can be done with privileged instruction:

```
change_keyboard_parameters:
    disable_interrupts
```

```
/* change things used by
    handle_keyboard_interrupt here */
```

enable_interrupts

exception implementation

detect condition (program error or external event)

save current value of PC somewhere

jump to exception handler (part of OS) jump done without program instruction to do so

exception implementation: notes

- I describe a simplified version
- real x86/x86-64 is a bit more complicated (mostly for historical reasons)

context

- all registers values %rax %rbx, ..., %rsp, ...
- condition codes
- program counter
- address space (map from program to real addresses)

context switch pseudocode

```
context switch(last, next):
  copy preexception pc last->pc
  mov rax,last->rax
  mov rcx, last->rcx
  mov rdx, last->rdx
  . . .
  mov next->rdx. rdx
  mov next->rcx, rcx
  mov next->rax, rax
  imp next->pc
```

applicatic	ons						
	standard library functions / shell commands						
	standard libraries and libc (C standard library) the					the shel	
	utility p	utility programs login					login
		system	call inter	face			
		kernel	CPU schee virtual me pipes	duler mory	filesystems device drivers swapping	netv sign 	vorking als
	hardwa	re inter	face				
hardware	m	emory m	anagement	unit	device controlle	rs .	

applications							
standard library functions / shell commands							
stan	dard libraries and libc (C standard library) the she						
utili	ty programs login login						
	system call interface						
user-mode hardware interface (limited)	CPU schedulerfilesystemsnetworkingkernelvirtual memorydevice driverssignalspipesswapping						
(initial)	kernel-mode hardware interface (complete)						
hardware	memory management unit device controllers						

applicatio	ns							
standard library functions / shell commands								
	standard libraries and libc (C standard library)						the shell	
	utility	progran	าร	login	I		login	
user-mode hardware interface (limited)		system	call inter	rface				
		kernel	CPU sche virtual me pipes	duler mory	filesystems device drivers swapping	netv sign 	working Ials	
	-)	kernel-mode hardware interface (complete)						
hardware	m	emory m	anagement	unit	device controlle	ers		

standa	standard library functions / shell commands						
standa	standard libraries and libc (C standard library) th						
utility	programs login	login					
user mode	system call interface						
hardware interface (limited)	CPU scheduler filesystems net kernel virtual memory device drivers sig pipes swapping	working nals					
(kernel-mode hardware interface (complete)						
hardware ⁿ	nemory management unit device controllers						

the OS?

applications							
standar	d librai	ry functior	ns / s	shell comman	ds		
standar	d librai	ries and	libc	(C standard libra	ary)	the shell	
utility p	progran	ns	login			login	
	system	call inter	face				1
user-mode hardware interface (limited)	kernel	CPU sched virtual mer pipes	luler mory	filesystems device drivers swapping	netv sigr 	working nals	the OS?
(kerne	el-mode ha	ardwa	are interface (com	nplete)	
hardware me	emory m	anagement	unit	device controlle	ers		47

aside: is the OS the kernel?

- OS = stuff that runs in kernel mode?
- OS = stuff that runs in kernel mode + libraries to use it?
- OS = stuff that runs in kernel mode + libraries + utility programs (e.g. shell, finder)?
- OS = everything that comes with machine?

no consensus on where the line is

each piece can be replaced separately...

exception implementation

detect condition (program error or external event)

save current value of PC somewhere

jump to exception handler (part of OS) jump done without program instruction to do so

exception implementation: notes

- I describe a simplified version
- real x86/x86-64 is a bit more complicated (mostly for historical reasons)

running the exception handler

hardware saves the old program counter (and maybe more)

identifies location of exception handler via table

then jumps to that location

OS code can save anything else it wants to , etc.