

Lecture 2 - System Programming Concepts

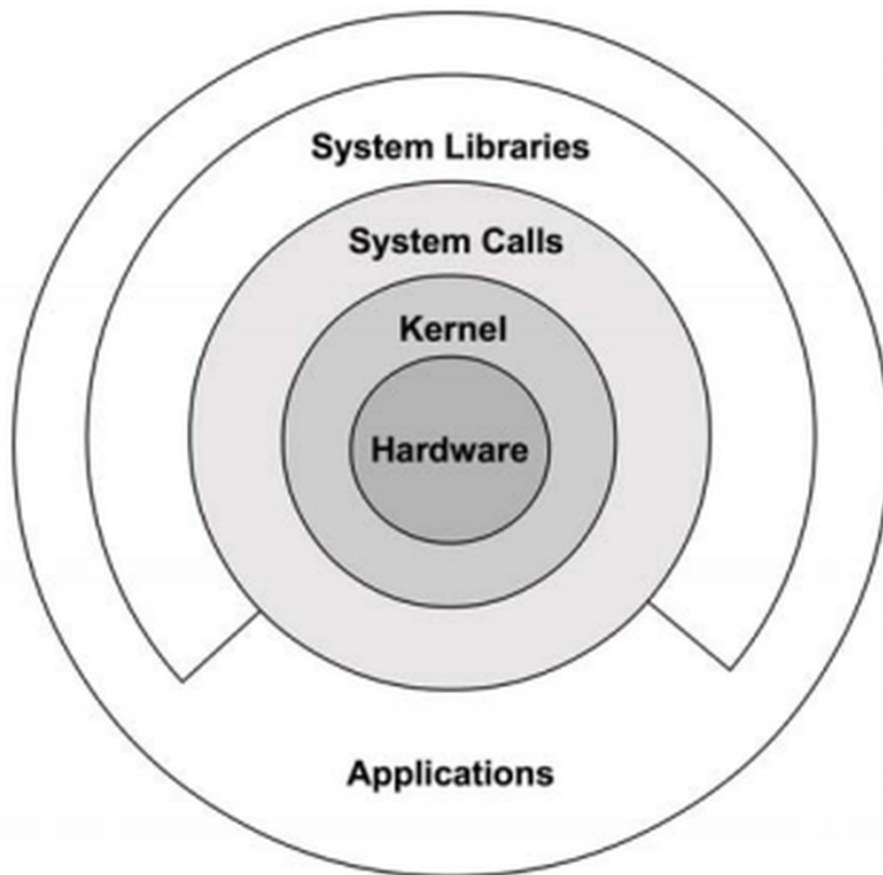
EE3233 Systems Programming for Engrs Reference: M. Kerrisk, The Linux Programming Interface



UTSA® The University of Texas at San Antonio™

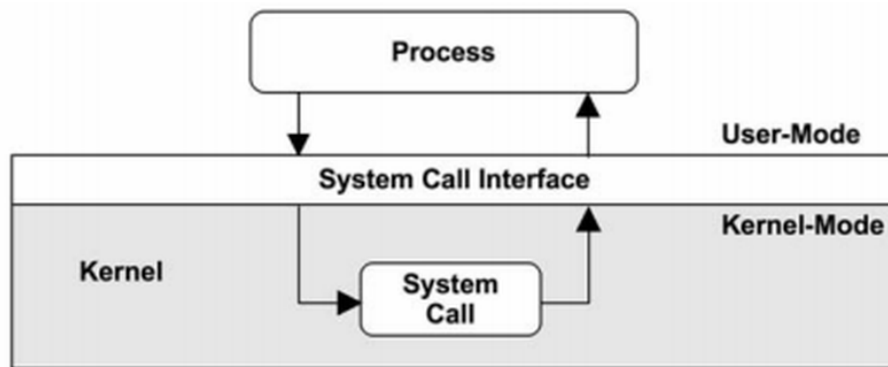
System Calls

- Entry point into kernel
- With System Calls, a process requests to kernel to perform some action on the process's behalf
 - Kernel makes a range of services (creating new process, performing I/O, creating a pipe for IPC) accessible to programs via the system call API



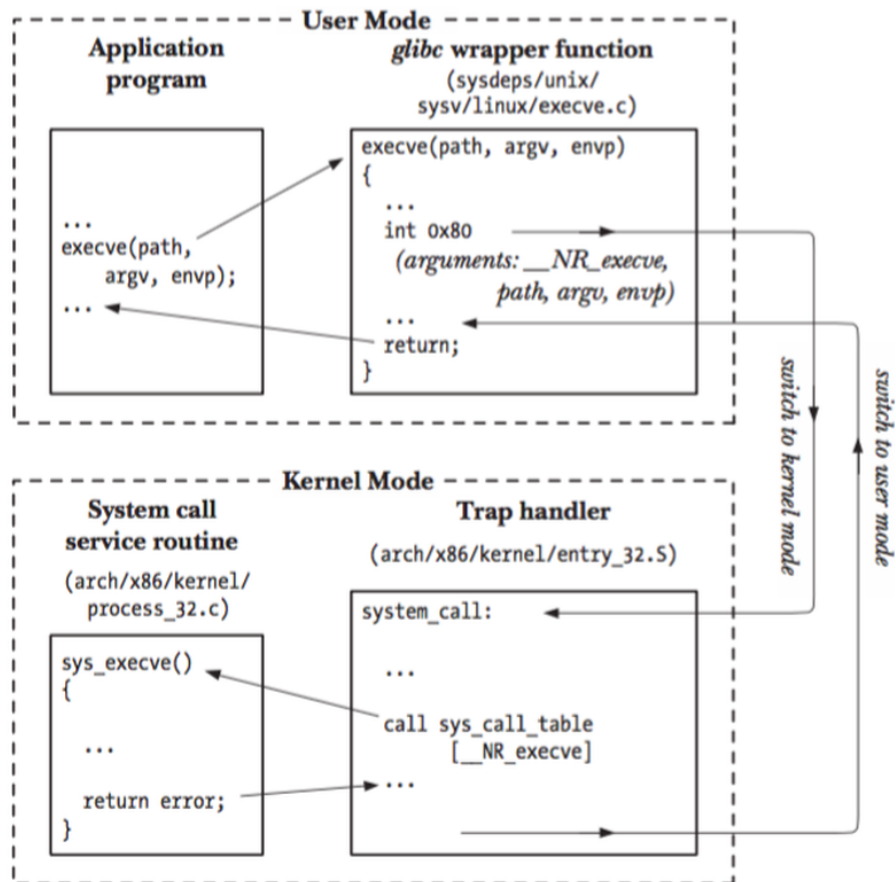
System Calls - Cont'd

- A system call changes processor state from user mode to kernel mode
 - so that CPU can access protected kernel memory
- The set of system calls is fixed
- Arguments of system call are transferred from user space to kernel space (vice versa)



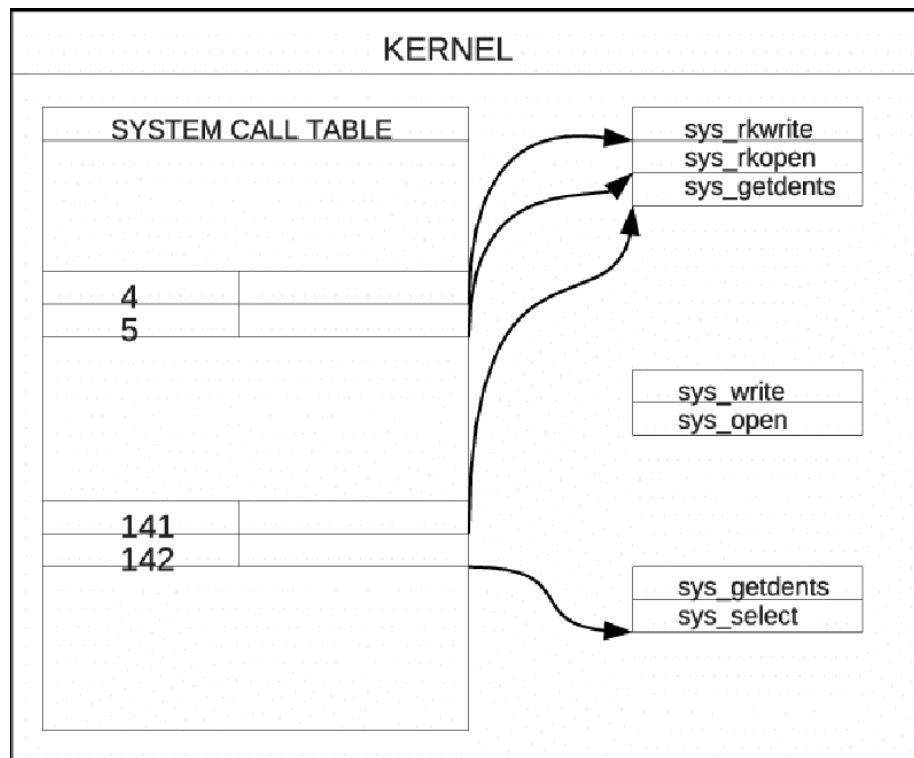
System Calls - Cont'd

1. Application makes a system call by invoking a wrapper function in the C library
2. Wrapper function makes all arguments available to trap-handling routine
 - These arguments are passed to wrapper via the stack
 - Wrapper copies arguments to register (kernel requires them in registers)
3. Wrapper copies system call number into a specific CPU register (%eax)
4. Wrapper executes a trap machine instruction (int 0x80) : switch user mode → kernel mode
5. Kernel invokes its `system__call()` routine located at "arch/x86/kernel/entry_32.S"
 - Saves register values onto the kernel stack
 - Checks the validity of the system call number
 - Invokes appropriate system call service routine
 - If the system call service routine has any arguments, it first checks their validity (addresses of pointers are valid)
6. Trap handler places the system call return value on the stack
 - If the system call service routine returned the error value, the wrapper function sets the global variable `errno` using this value
 - Simultaneously returning to user mode



System Calls - Cont'd

`execve()` : system call number 11 (`__NR_execve`) * in `system_call_table` vector, entry 11 contains the address of `sys_execve()`, service routine for this system call (`sys_xxx()` is a typical name in Linux) * “invoking the system call `xyz()`” * means “calling the wrapper function that invokes the system call `xyz()`”



Library Functions

- Functions that constitute standard C library
- Many library functions don't use system calls
 - e.g., string manipulation functions
- Some library functions are layered on top of system calls
 - e.g., `fopen()` library function uses the `open()` system call
 - designed to provide caller friendly interface
 - * e.g., `printf()` provides output formatting while `write()` just outputs a block of bytes
 - * e.g., `malloc()` and `free()` performs various bookkeeping tasks while `brk()` system call does not

Standard C Library (glibc)

- The most commonly used implementation on Linux is GNU C library (glibc)
- Version of glibc

```
$ /lib/x86_64-linux-gnu/libc.so.6

root@90dca184ddd2:/notebooks/temp# /lib/x86_64-linux-gnu/libc.so.6
GNU C Library (Ubuntu GLIBC 2.23-0ubuntu10) stable release version 2.23, by Roland McGrath et al.
Copyright (C) 2016 Free Software Foundation, Inc.
This is free software; see the source for copying conditions.
There is NO warranty; not even for MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.
Compiled by GNU CC version 5.4.0 20160609.
Available extensions:
crypt add-on version 2.1 by Michael Glad and others
GNU Libidn by Simon Josefsson
Native POSIX Threads Library by Ulrich Drepper et al
BIND-8.2.3-T5B
libc ABIs: UNIQUE IFUNC
For bug reporting instructions, please see:
<https://bugs.launchpad.net/ubuntu/+source/glibc/+bugs>.
```

Handling Errors from System Calls and Library Functions

- They return status value indicating success or fail
 - This status value should always be checked
 - Excluding this check to save time is a false economy
-

Handling system call errors

Usually, an error is indicated by a return of -1

```
fd = open(pathname, flags, mode);  /* system call to open a file */
if ( fd == -1 )
{
    /* Code to handle the error */
}
***
if ( close(fd) == -1 )
{
    /* Code to handle the error */
}
```

Handling system call errors

- When a system call fails, it sets the global variable *errno* to a positive value (specific error)

– you should include `<errno.h>` header (declaration of *errno*)

```
cnt = read(fd, buf, numbytes);
if (cnt == -1)
{
    if (errno == EINTR) {
        fprintf(stderr, "read was interrupted by a signal\n");
    }
    else {
        /* Some other error occurred */
    }
}

/usr/include/asm-generic/errno-base.h
```

Handling system call errors

- Common action after a failed system call is to print an error message based on the *errno* value : use `perror()`
- `perror()` prints the string pointed by `msg` and message corresponding to the current value of *errno*

Prototype:

```
#include <stdio.h>

void perror (const char *msg);

fd = open (pathname, flags, mode);
if ( fd == -1 )
{
    perror ("open");
    exit (EXIT_FAILURE);
}
```

System Data Type

- Most of the standard system data types have names ending in `**_t**`

```
typedef int pid_t;          /* intended to represent process IDs */

pid_t mypid;
```

System Data Type

| Data type | SUSv3 type requirement | Description |
|-------------------|---|--|
| <i>blkcnt_t</i> | signed integer | File block count (Section 15.1) |
| <i>blksize_t</i> | signed integer | File block size (Section 15.1) |
| <i>cc_t</i> | unsigned integer | Terminal special character (Section 62.4) |
| <i>clock_t</i> | integer or real-floating | System time in clock ticks (Section 10.7) |
| <i>clockid_t</i> | an arithmetic type | Clock identifier for POSIX.1b clock and timer functions (Section 23.6) |
| <i>comp_t</i> | not in SUSv3 | Compressed clock ticks (Section 28.1) |
| <i>dev_t</i> | an arithmetic type | Device number, consisting of major and minor numbers (Section 15.1) |
| <i>DIR</i> | no type requirement | Directory stream (Section 18.8) |
| <i>fd_set</i> | structure type | File descriptor set for <i>select()</i> (Section 63.2.1) |
| <i>fsblkcnt_t</i> | unsigned integer | File-system block count (Section 14.11) |
| <i>fsfilcnt_t</i> | unsigned integer | File count (Section 14.11) |
| <i>gid_t</i> | integer | Numeric group identifier (Section 8.3) |
| <i>id_t</i> | integer | A generic type for holding identifiers; large enough to hold at least <i>pid_t</i> , <i>uid_t</i> , and <i>gid_t</i> |
| <i>in_addr_t</i> | 32-bit unsigned integer | IPv4 address (Section 59.4) |
| <i>in_port_t</i> | 16-bit unsigned integer | IP port number (Section 59.4) |
| <i>ino_t</i> | unsigned integer | File i-node number (Section 15.1) |
| <i>key_t</i> | an arithmetic type | System V IPC key (Section 45.2) |
| <i>mode_t</i> | integer | File permissions and type (Section 15.1) |
| <i>mqd_t</i> | no type requirement, but shall not be an array type | POSIX message queue descriptor |
| <i>msglen_t</i> | unsigned integer | Number of bytes allowed in System V message queue (Section 46.4) |

| Data type | SUSv3 type requirement | Description |
|---------------------|---|---|
| <i>msgqnum_t</i> | unsigned integer | Counts of messages in System V message queue (Section 46.4) |
| <i>nfds_t</i> | unsigned integer | Number of file descriptors for <i>poll()</i> (Section 63.2.2) |
| <i>nlink_t</i> | integer | Count of (hard) links to a file (Section 15.1) |
| <i>off_t</i> | signed integer | File offset or size (Sections 4.7 and 15.1) |
| <i>pid_t</i> | signed integer | Process ID, process group ID, or session ID (Sections 6.2, 34.2, and 34.3) |
| <i>ptrdiff_t</i> | signed integer | Difference between two pointer values, as a signed integer |
| <i>rlim_t</i> | unsigned integer | Resource limit (Section 36.2) |
| <i>sa_family_t</i> | unsigned integer | Socket address family (Section 56.4) |
| <i>shmatt_t</i> | unsigned integer | Count of attached processes for a System V shared memory segment (Section 48.8) |
| <i>sig_atomic_t</i> | integer | Data type that can be atomically accessed (Section 21.1.3) |
| <i>siginfo_t</i> | structure type | Information about the origin of a signal (Section 21.4) |
| <i>sigset_t</i> | integer or structure type | Signal set (Section 20.9) |
| <i>size_t</i> | unsigned integer | Size of an object in bytes |
| <i>socklen_t</i> | integer type of at least 32 bits | Size of a socket address structure in bytes (Section 56.3) |
| <i>speed_t</i> | unsigned integer | Terminal line speed (Section 62.7) |
| <i>ssize_t</i> | signed integer | Byte count or (negative) error indication |
| <i>stack_t</i> | structure type | Description of an alternate signal stack (Section 21.3) |
| <i>suseconds_t</i> | signed integer allowing range [-1, 1000000] | Microsecond time interval (Section 10.1) |
| <i>tcflag_t</i> | unsigned integer | Terminal mode flag bit mask (Section 62.2) |
| <i>time_t</i> | integer or real-floating | Calendar time in seconds since the Epoch (Section 10.1) |
| <i>timer_t</i> | an arithmetic type | Timer identifier for POSIX.1b interval timer functions (Section 23.6) |
| <i>uid_t</i> | integer | Numeric user identifier (Section 8.1) |