

# **Linux/UNIX** **System Programming** **Essentials**

**Michael Kerrisk**

**[man7.org](http://man7.org)**

**November 2021**



**©2021, man7.org Training and Consulting /  
Michael Kerrisk. All rights reserved.**

These training materials have been made available for personal, noncommercial use. Except for personal use, no part of these training materials may be printed, reproduced, or stored in a retrieval system. These training materials may not be redistributed by any means, electronic, mechanical, or otherwise, without prior written permission of the author. These training materials may not be used to provide training to others without prior written permission of the author.

Every effort has been made to ensure that the material contained herein is correct, including the development and testing of the example programs. However, no warranty is expressed or implied, and the author shall not be liable for loss or damage arising from the use of these programs. The programs are made available under Free Software licenses; see the header comments of individual source files for details.

For information about this course, visit  
<http://man7.org/training/>.

For inquiries regarding training courses, please contact us at  
[training@man7.org](mailto:training@man7.org).

Please send corrections and suggestions for improvements to this course material to [training@man7.org](mailto:training@man7.org).

For information about *The Linux Programming Interface*, please visit <http://man7.org/tlpi/>.



## Short table of contents

---

1	Course Introduction	1-1
2	Fundamental Concepts	2-1
3	File I/O and Files	3-1
4	Processes	4-1
5	Signals	5-1
6	Process Lifecycle	6-1
7	System Call Tracing with strace	7-1



## Detailed table of contents

---

1	Course Introduction	1-1
1.1	Course overview	1-3
1.2	Course materials and resources	1-9
1.3	Introductions	1-14
2	Fundamental Concepts	2-1
2.1	System calls and library functions	2-3
2.2	Error handling	2-10
2.3	System data types	2-17
2.4	Notes on code examples	2-22
3	File I/O and Files	3-1
3.1	File I/O overview	3-3
3.2	open(), read(), write(), and close()	3-8
3.3	Retrieving file information: stat()	3-23
4	Processes	4-1
4.1	Process IDs	4-3
4.2	Process memory layout	4-6
4.3	Command-line arguments	4-9
4.4	The environment list	4-11

## Detailed table of contents

---

4.5	The /proc filesystem	4-16
5	Signals	5-1
5.1	Overview of signals	5-3
5.2	Signal dispositions	5-8
5.3	Useful signal-related functions	5-16
5.4	Signal handlers	5-20
5.5	Signal sets, the signal mask, and pending signals	5-25
5.6	Designing signal handlers	5-37
6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: fork()	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The SIGCHLD signal	6-38
6.7	Executing programs: execve()	6-42
7	System Call Tracing with strace	7-1

## Detailed table of contents

---

7.1	Getting started	7-3
7.2	Tracing child processes	7-10
7.3	Filtering strace output	7-14
7.4	System call tampering	7-20
7.5	Further strace options	7-26

# Linux System Programming Essentials

## Course Introduction

Michael Kerrisk, man7.org © 2021

mtk@man7.org

November 2021

### Outline

---

1	Course Introduction	1-1
1.1	Course overview	1-3
1.2	Course materials and resources	1-9
1.3	Introductions	1-14



# Outline

---

1	Course Introduction	1-1
1.1	Course overview	1-3
1.2	Course materials and resources	1-9
1.3	Introductions	1-14

## Course prerequisites

---

- Prerequisites
  - (Good) reading knowledge of C
  - Can log in to Linux / UNIX and use basic commands
- Knowledge of *make(1)* is helpful
  - (Can do a short tutorial during first practical session for those new to *make*)
- Assumptions
  - You are familiar with commonly used parts of standard C library
    - e.g., *stdio* and *malloc* packages

## Course goals

---

- Aimed at programmers building/understanding low-level applications
- Gain strong understanding of programming API that kernel presents to user-space
  - System calls
  - Relevant C library functions
  - Other interfaces (e.g., /proc)
  - Necessarily, we sometimes delve into inner workings of kernel
    - (But... not an internals course)
- Course topics
  - Course flyer
  - For more detail, see TOC in course books

## Lab sessions

---

- Lots of lab sessions...
- **Pair/group work is strongly encouraged!**
  - Pairs typically get through practical sessions faster
    - ⇒ we will go faster as a group, and cover more topics
- **Read each exercise thoroughly** before starting
  - Past experience has shown me the traps that people often fall into with various exercises
  - ⇒ exercise descriptions often include **important hints**
- Lab sessions are **not** instructor down time...
  - ⇒ One-on-one questions about course material or exercises
- Looking for homework?
  - ⇒ Chapters usually have additional exercises

## Coding exercises

---

- For coding exercises, you can use any suitable programming language in which you are proficient
  - C/C++ (easiest...)
  - Go, D, Rust, & other languages that compile to native machine code
  - Most features can also be exercised from scripting languages such as Python, Ruby, and Perl
- **Template** solutions are provided for most coding exercises
  - Filenames: `ex.*.c`
  - Look for “FIXMEs” to see what parts you must complete
  - ⚠ You will need to edit the corresponding Makefile to add a new target for the executable
- **Solutions** will be mailed out shortly after end of course

## Lab sessions: some thoughts on building code

---

- Many warnings indicate real problems with your code; fix them
  - And the “harmless errors” create noise that hides the serious warnings; fix them
  - **This is a good thing: `cc -Werror`**
    - Treat all warnings as errors
- Rather than writing lots of code before first compile, use a frequent edit-save-build cycle to catch compiler errors early
  - Try running the following in a separate window as you edit:

```
$ while inotifywait -q . ; do echo -e '\n\n'; make; done
```
  - *inotifywait* is provided in the *inotify-tools* package
  - (The *echo* command just injects some white space between each build)

# Outline

---

1	Course Introduction	1-1
1.1	Course overview	1-3
1.2	Course materials and resources	1-9
1.3	Introductions	1-14

## Course materials

---

- Source code tarball
  - Location sent by email
  - Unpacked source code is a Git repository; you can commit/revert changes, etc.
- Slides / course book
- Kerrisk, M.T. 2010. *The Linux Programming Interface* (TLPI), No Starch Press.
  - Slides frequently reference TLPI in bottom RHS corner
  - Further info on TLPI: <http://man7.org/tlpi/>
    - API changes since publication:  
[http://man7.org/tlpi/api\\_changes/](http://man7.org/tlpi/api_changes/)

## Other resources

---

- POSIX.1-2001 / SUSv3: <http://www.unix.org/version3/>
- POSIX.1-2008 / SUSv4: <http://www.unix.org/version4/>
- Man pages
  - Section 2: system calls
  - Section 3: library functions
  - Latest version online at <http://man7.org/linux/man-pages/>
  - Latest tarball downloadable at <https://www.kernel.org/doc/man-pages/download.html>

## Books

---

- General:
  - Stevens, W.R., and Rago, S.A. 2013. *Advanced Programming in the UNIX Environment (3rd edition)*. Addison-Wesley.
    - <http://www.apuebook.com/>
- POSIX threads:
  - Butenhof, D.R. 1996. *Programming with POSIX Threads*. Addison-Wesley.
- TCP/IP and network programming:
  - Fall, K.R. and Stevens, W.R. 2013. *TCP/IP Illustrated, Volume 1: The Protocols (2nd Edition)*. Addison-Wesley.
  - Stevens, W.R., Fenner, B., and Rudoff, A.M. 2004. *UNIX Network Programming, Volume 1 (3rd edition): The Sockets Networking API*. Addison-Wesley.
    - <http://www.unpbook.com/>
  - Stevens, W.R. 1999. *UNIX Network Programming, Volume 2 (2nd edition): Interprocess Communications*. Prentice Hall.
    - <http://www.kohala.com/start/unpv22e/unpv22e.html>
- Operating systems:
  - Tanenbaum, A.S., and Woodhull, A.S. 2006. *Operating Systems: Design And Implementation (3rd edition)*. Prentice Hall.
    - (The Minix book)
  - Comer, D. 2015. *Operating System Design: The Xinu Approach (2nd edition)*

# Common abbreviations used in slides

---

The following abbreviations are sometimes used in the slides:

- ACL: access control list
- COW: copy-on-write
- CV: condition variable
- CWD: current working directory
- EA: extended attribute
- EOF: end of file
- FD: file descriptor
- FS: filesystem
- FTM: feature test macro
- GID: group ID
  - rGID, eGID, sGID, fsGID
- iff: “if and only if”
- IPC: interprocess communication
- KSE: kernel scheduling entity
- MQ: message queue
- MQD: message queue descriptor
- NS: namespace
- OFD: open file description
- PG: process group
- PID: process ID
- PPID: parent process ID
- SHM: shared memory
- SID: session ID
- SEM: semaphore
- SUS: Single UNIX specification
- UID: user ID
  - rUID, eUID, sUID, fsUID

## Outline

---

1	Course Introduction	1-1
1.1	Course overview	1-3
1.2	Course materials and resources	1-9
1.3	Introductions	1-14

## Introductions: me

---

- Programmer, trainer, writer
- UNIX since 1987, Linux since mid-1990s
- Active contributor to Linux
  - API review, testing, and documentation
    - API design and design review
    - Lots of testing, lots of bug reports, a few kernel patches
  - Maintainer of Linux *man-pages* project
    - Documents kernel-user-space + C library APIs
    - Contributor since 2000
    - As maintainer:  $\approx 22k$  commits, 194 releases since 2004
    - Author/coauthor of  $\approx 440$  out of  $\approx 1060$  man pages
- Kiwi in .de
  - (mtk@man7.org, PGP: 4096R/3A35CE5E)
  - @mkerrisk (feel free to tweet about the course as we go...)
  - <http://linkedin.com/in/mkerrisk>

## Introductions: you

---

In brief:

- Who are you?
  - If virtual: where are you?
- (Optionally:) any special goals for the course?
- Two things you like to do when you are not in front of a keyboard, and one thing you don't like doing...

## Questions policy

- General policy: ask questions any time
  - Either on Slack channel
  - Or out loud
    - (But, wait for a quiet point)
- In the event that questions slow us down too much, I may say: “batch your questions until next *Question penguin* slide”

## Notes

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and extend across the width of the page. There are no margins, text, or other markings on the paper.



Notes

Notes

# Linux System Programming Essentials

# Fundamental Concepts

Michael Kerrisk, man7.org © 2021

mtk@man7.org

November 2021

## Outline

---

2	Fundamental Concepts	2-1
2.1	System calls and library functions	2-3
2.2	Error handling	2-10
2.3	System data types	2-17
2.4	Notes on code examples	2-22

# Outline

---

2	Fundamental Concepts	2-1
2.1	System calls and library functions	2-3
2.2	Error handling	2-10
2.3	System data types	2-17
2.4	Notes on code examples	2-22

## System calls

---

System call == controlled entry point into kernel code

- Request to kernel to perform some task on caller's behalf
- *syscalls(2)* man page lists (nearly) all system calls
- Documented in Section 2 of man pages (notation: *stat(2)*)

# Steps in the execution of a system call

- ① Program calls wrapper function in C library
- ② Wrapper function packages syscall arguments into hardware registers
- ③ Wrapper function puts syscall number into a register
  - Each syscall has a unique number
- ④ Wrapper function traps to kernel mode
  - e.g., `syscall` instruction on x86-64 (or `sysenter` for 32-bit)
- ⑤ Kernel then executes syscall handler:
  - Checks validity of syscall number
  - Invokes **service routine** corresponding to syscall number
    - Checks arguments, **does real work**, returns a result status
  - Places syscall return value in a register
  - Switches back to user mode, passing control back to wrapper
    - E.g., `sysret` instruction on x86-64
- ⑥ Wrapper function examines syscall return value; on error, copies return value to `errno`

## System calls are expensive!

$10^9$  calls to...

simple user-space function returning *int*  $\Rightarrow$  1.5 seconds

*getppid()* system call  $\Rightarrow$  41 340 seconds

(The page table isolation patches to mitigate Spectre, Meltdown, etc. have resulted in a big performance hit on system calls)

(*getppid()*, which returns process ID of caller's parent, is one of the simplest system calls)

(Linux 5.4, x86-64; Intel Core i7-8850H; `progconc/syscall_speed.c`)

## Library functions

---

- Library function == one of multitude of functions in Standard C Library
- Diverse range of tasks:
  - I/O
  - Dynamic memory allocation
  - Math
  - String processing
  - etc.
- Documented in Section 3 of man pages (notation: *fopen(3)*)
- Some library functions employ system calls
- Many library functions make no use of system calls

[TLPI §3.2]

## The C library

---

- Each C environment has its own implementation of standard C library
- Linux has multiple implementations
- **GNU C library (glibc)** is most widely used
  - Full implementation of POSIX APIs, plus many extensions
  - <http://www.gnu.org/software/libc/>

[TLPI §3.3]

# The C library

---

- Other Linux C libraries target embedded platforms or the creation of small binaries:
  - musl (“mussel”) libc (<http://www.musl-libc.org/>)
    - Under active development (release 1.0 in 2014)
    - [http://wiki.musl-libc.org/wiki/Functional\\_differences\\_from\\_glibc](http://wiki.musl-libc.org/wiki/Functional_differences_from_glibc)
  - uclibc (<http://www.uclibc.org/>) [inactive?]
  - dietlibc (<http://www.fefe.de/dietlibc/>) [inactive?]
  - A comparison: [http://www.etalabs.net/compare\\_libcs.html](http://www.etalabs.net/compare_libcs.html)
- (C library on Android is Bionic)
- We’ll presume the use of glibc

## Outline

---

2	Fundamental Concepts	2-1
2.1	System calls and library functions	2-3
2.2	Error handling	2-10
2.3	System data types	2-17
2.4	Notes on code examples	2-22

## Error handling

---

- Most system calls and library functions return a status indicating success or failure
- In failure, most system calls:
  - Return `-1`
  - Place integer value in global variable `errno` to indicate cause
- Some library functions follow same convention
- Often, we'll omit return values from slides, where they follow usual conventions
  - Check man pages for details

## Error handling

---

- Return status should **always** be tested
- ⚠ Inspect `errno` only if result status indicates failure
  - APIs do not reset `errno` to 0 on success
  - A successful call may modify `errno` (POSIX allows this)
  - E.g., this is wrong:

```
fd = open(pathname, O_RDONLY);  
printf("open() has returned\n");    // Might modify errno!  
  
if (fd == -1) {  
    perror("open");  
    exit(EXIT_FAILURE);  
}
```

- When an API call fails, *errno* is set to indicate cause
- Integer value, global variable
  - In multithreading environment, each thread has private *errno*
- Error numbers in *errno* are  $> 0$
- `<errno.h>` defines symbolic names for error numbers

```
#define EPERM      1  /* Operation not permitted */
#define ENOENT     2  /* No such file or directory */
#define ESRCH      3  /* No such process */
...
```

- *errno(1)* command can be used to search for errors by number, name, or substring in textual message
  - Part of *moreutils* package

## Checking for errors

```
1 cnt = read(fd, buf, numbytes);
2
3 if (cnt == -1) {      /* Was there an error? */
4     if (errno == EINTR)
5         fprintf(stderr,
6             "read() was interrupted by a signal\n");
7     else if (errno == EBADF)
8         fprintf(stderr,
9             "read() given bad file descriptor\n");
10    else {
11        /* Some other error occurred */
12    }
13 }
```



## Displaying error messages

---

```
#include <stdio.h>
void perror(const char *msg);
```

- Outputs to *stderr*:
  - *msg* + ":" + string corresponding to value in *errno*
  - E.g., if *errno* contains EBADF, *perror("close")* would display:  
close: Bad file descriptor
- Simple error handling:

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

## Displaying error messages

---

```
#include <string.h>
char *strerror(int errnum);
```

- Returns an error string corresponding to error in *errnum*
  - Same string as printed by *perror()*
- Unknown error number?  $\Rightarrow$  "*Unknown error nnn*"
  - Or NULL on some systems

# Outline

---

2	Fundamental Concepts	2-1
2.1	System calls and library functions	2-3
2.2	Error handling	2-10
2.3	System data types	2-17
2.4	Notes on code examples	2-22

## System data types

---

- Various system info needs to be represented in C
  - Process IDs, user IDs, file offsets, etc.
- Using native C data types (e.g., *int*, *long*) in application code would be nonportable; e.g.:
  - `sizeof(long)` might be 4 on one system, but 8 on another
  - One system might use *int* for PIDs, while another uses *long*
  - Even on same system, things may change across versions
    - E.g., in kernel 2.4, Linux switched from 16 to 32-bit UIDs
- ⇒ POSIX defines system data types:
  - Implementations must suitably define each system data type
  - Defined via `typedef`; e.g., `typedef int pid_t`
    - Most types have names suffixed “\_t”
  - Applications should use these types; e.g., `pid_t mypid`;
    - ⇒ will compile to correct types on any conformant system

[TLPI §3.6.2]

## Examples of system data types

Data type	POSIX type requirement	Description
<i>uid_t</i>	Integer	User ID
<i>gid_t</i>	Integer	Group ID
<i>pid_t</i>	Signed integer	Process ID
<i>id_t</i>	Integer	Generic ID type; can hold <i>pid_t</i> , <i>uid_t</i> , <i>gid_t</i>
<i>off_t</i>	Signed integer	File offset or size
<i>sigset_t</i>	Integer or structure	Signal set
<i>size_t</i>	Unsigned integer	Size of object (in bytes)
<i>ssize_t</i>	Signed integer	Size of object or error indication
<i>time_t</i>	Integer/real-floating	Time in seconds since Epoch
<i>timer_t</i>	Arithmetic type	POSIX timer ID

(Arithmetic type  $\in$  integer or floating type)

## Printing system data types

- Need to take care when passing system data types to *printf()*
- Example: *pid\_t* can be *short*, *int*, or *long*
- Suppose we write:

```
printf("My PID is: %d\n", getpid());
```

- Works fine if:
  - *pid\_t* is *int*
  - *pid\_t* is *short* (C promotes *short* argument to *int*)
- But **what if *pid\_t* is *long*** (and *long* is bigger than *int*)?
  - $\Rightarrow$  argument exceeds range understood by format specifier (top bytes will be lost)

## Printing system data types

---

- On virtually all implementations, most integer system data types are *long* or smaller

- ⇒ Promote to *long* when printing system data types

```
printf("My PID is: %ld\n", (long) getpid());
```

- Most notable exception: *off\_t* is typically *long long*

- Promote to *long long* for *printf()*

```
printf("Offset is %lld\n",  
      (long long) lseek(fd, 0, SEEK_CUR));
```

- Can also use *%zu* and *%zd* for *size\_t* and *ssize\_t*
- C99 has *intmax\_t* (*uintmax\_t*) with *%jd* (*%ju*) *printf()* specifier
  - Solution for all integer types, but not on pre-C99 systems
  - Must include `<stdint.h>` to get these type definitions

## Outline

---

2	Fundamental Concepts	2-1
2.1	System calls and library functions	2-3
2.2	Error handling	2-10
2.3	System data types	2-17
2.4	Notes on code examples	2-22

## Code examples presented in course

---

- **Code tarball** == code from TLPI + further code for course
- **Examples on slides edited/excerpted** for brevity
  - E.g., error-handling code may be omitted
- Slides always show **pathname for full source code**
  - Full source code always includes error-handling code
- Code license:
  - GNU GPL v3 for programs
  - GNU Lesser GPL v3 for library functions
  - <http://www.gnu.org/licenses/#GPL>
    - Understanding Open Source and Free Software Licensing; A.M. St Laurent, 2004
    - Open Source Licensing: Software Freedom and Intellectual Property Law; L. Rosen, 2004
    - Open Source Software: Rechtliche Rahmenbedingungen der Freien Software; Till Jaeger, 2020

## Example code lib/ subdirectory

---

- lib/ subdirectory contains code of a few functions commonly used in examples
- *camel/Case* function name?
  - ⇒ It's mine

## Common header file

---

- Many code examples make use of header file `tlpi_hdr.h`
- Goal: make code examples a little shorter
- `tlpi_hdr.h`:
  - Includes a few frequently used header files
  - Includes declarations of some error-handling functions

[TLPI §3.5.2]

## Error-handling functions used in examples

---

- Could handle errors as follows:

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

- To save some effort, I define some simple error-handling functions

## Error-handling functions used in examples

```
#include "tldpi_hdr.h"
errExit(const char *format, ...);
```

- Prints error message on *stderr* that includes:
  - Symbolic name for *errno* value (via some trickery)
  - *strerror()* description for current *errno* value
  - Text from the *printf()*-style message supplied in arguments
  - A terminating newline
- Terminates program with exit status EXIT\_FAILURE (1)
- Example:

```
if (close(fd) == -1)
    errExit("close (fd=%d)", fd);
```

might produce:

```
ERROR [EBADF Bad file descriptor] close (fd=5)
```

## Error-handling functions used in examples

```
#include "tldpi_hdr.h"
errMsg(const char *format, ...);
```

- Like *errExit()*, but does not terminate program

```
#include "tldpi_hdr.h"
fatal(const char *format, ...);
```

- Displays a *printf()*-style message + newline
- Terminates program with exit status EXIT\_FAILURE (1)

## Building the sample code

- You can manually compile the example programs, but there is also a **Makefile** in each directory
- ⇒ Typing `make` in source code root directory builds all programs in all subdirectories
- If you encounter build errors relating to ACLs, capabilities, or SELinux, see <http://man7.org/tlpi/code/faq.html>
  - Preferred solution is to install the necessary packages:
    - Debian derivatives: *libcap-dev*, *libacl1-dev*, *libselinux1-dev*
    - RPM-based systems: *libcap-devel*, *libacl-devel*, *libselinux-devel*
  - If you can't install these packages, then:

```
cd lib
sh Build_lib.sh      # Ignore any errors you see
```

and then do *make* in individual directories as needed

## Using library functions from the sample code

To use my library functions in your code:

- **Include** `tlpi_hdr.h` in your C source file
  - Located in `lib/` subdirectory in source code
- **Link against my library**, `libtlpi.a`, located in source code root directory
  - To build library, run `make` in the source code root directory or in `lib/` subdirectory
- **Method 1:** Place your program in one of “my” directories, add target to corresponding Makefile, and build using `make`
- **Method 2:** Manually compile with the following command:

```
cc -Isrc-root/lib yourprog.c src-root/libtlpi.a
```

- *src-root* must be replaced with the absolute or relative path of source code root directory



Notes

Notes

Linux System Programming Essentials

# File I/O and Files

Michael Kerrisk, man7.org © 2021

mtk@man7.org

November 2021

## Outline

---

3	File I/O and Files	3-1
3.1	File I/O overview	3-3
3.2	open(), read(), write(), and close()	3-8
3.3	Retrieving file information: stat()	3-23

## Outline

---

3	File I/O and Files	3-1
3.1	File I/O overview	3-3
3.2	open(), read(), write(), and close()	3-8
3.3	Retrieving file information: stat()	3-23

## Files

---

- “On UNIX, everything is a file”
  - More correctly: “everything is a file descriptor”
- Note: the term **file** can be ambiguous:
  - A **generic term**, covering disk files, directories, sockets, FIFOs, terminals and other devices and so on
  - Or specifically, a **disk file** in a filesystem
  - To clearly distinguish the latter, the term **regular file** is sometimes used

## System calls versus *stdio*

- C programs usually use *stdio* package for file I/O
- Library functions layered on top of I/O system calls

System calls	Library functions
file descriptor ( <i>int</i> )	file stream ( <i>FILE *</i> )
<i>open()</i> , <i>close()</i>	<i>fopen()</i> , <i>fclose()</i>
<i>lseek()</i>	<i>fseek()</i> , <i>ftell()</i>
<i>read()</i>	<i>fgets()</i> , <i>fscanf()</i> , <i>fread()</i> ...
<i>write()</i>	<i>fputs()</i> , <i>fprintf()</i> , <i>fwrite()</i> , ...
—	<i>feof()</i> , <i>ferror()</i>

- We presume understanding of *stdio*;  $\Rightarrow$  focus on system calls

## File descriptors

- All I/O is done using file descriptors (FDs)
  - nonnegative integer that identifies an open file
- Used for all types of files
  - terminals, regular files, pipes, FIFOs, devices, sockets, ...
- 3 FDs are normally available to programs run from shell:
  - (POSIX names are defined in `<unistd.h>`)

FD	Purpose	POSIX name	<i>stdio</i> stream
0	Standard input	STDIN_FILENO	<i>stdin</i>
1	Standard output	STDOUT_FILENO	<i>stdout</i>
2	Standard error	STDERR_FILENO	<i>stderr</i>

# Key file I/O system calls

---

Four fundamental calls:

- *open()*: open a file, optionally creating it if needed
  - Returns file descriptor used by remaining calls
- *read()*: input
- *write()*: output
- *close()*: close file descriptor

## Outline

---

3	File I/O and Files	3-1
3.1	File I/O overview	3-3
3.2	<i>open()</i> , <i>read()</i> , <i>write()</i> , and <i>close()</i>	3-8
3.3	Retrieving file information: <i>stat()</i>	3-23

## *open()*: opening a file

---

```
#include <sys/stat.h>
#include <fcntl.h>
int open(const char *pathname, int flags,
        ... /* mode_t mode */);
```

- Opens existing file / creates and opens new file
- Arguments:
  - *pathname* identifies file to open
  - *flags* controls semantics of call
    - e.g., open an existing file vs create a new file
  - *mode* specifies permissions when creating new file
- Returns: a file descriptor (nonnegative integer)
  - (Guaranteed to be lowest available FD)

[TLPI §4.3]

## *open()* flags argument

---

Created by ORing (|) together:

- Access mode
  - Specify exactly one of `O_RDONLY`, `O_WRONLY`, or `O_RDWR`
- File creation flags (bit flags)
- File status flags (bit flags)

[TLPI §4.3.1]

## File creation flags

---

- **File creation flags:**
  - Affect behavior of *open()* call
  - Can't be retrieved or changed
- Examples:
  - `O_CREAT`: create file if it doesn't exist
    - *mode* argument must be specified
    - Without `O_CREAT`, can open only an existing file (else: `ENOENT`)
  - `O_EXCL`: create "exclusively"
    - Give an error (`EEXIST`) if file already exists
    - Only meaningful with `O_CREAT`
  - `O_TRUNC`: truncate existing file to zero length

## File status flags

---

- **File status flags:**
  - Affect semantics of subsequent file I/O
  - Can be retrieved and modified using *fcntl()*
- Examples:
  - `O_APPEND`: always append writes to end of file
  - `O_SYNC`: make file writes synchronous
  - `O_NONBLOCK`: nonblocking I/O

## *open()* examples

- Open existing file for reading:

```
fd = open("script.txt", O_RDONLY);
```

- Open file for read-write, create if necessary, ensure we are creator:

```
fd = open("myfile.txt", O_CREAT | O_EXCL | O_RDWR,  
          S_IRUSR | S_IWUSR); /* rw----- */
```

- Open file for writing, creating if necessary:

```
fd = open("myfile.txt", O_CREAT | O_WRONLY, S_IRUSR);
```

- File opened for writing, but created with only read permission!
- Open for writing, create if necessary, truncate, always append writes:

```
fd = open("log", O_CREAT | O_WRONLY | O_TRUNC | O_APPEND,  
          S_IRUSR | S_IWUSR);
```

- (`O_TRUNC` plus `O_APPEND` could be useful if another process is also doing writes at the end of the file)

## *read()*: reading from a file

```
#include <unistd.h>  
ssize_t read(int fd, void *buffer, size_t count);
```

- Arguments:
  - *fd*: file descriptor
  - *buffer*: pointer to buffer to store data
    - ⚠ No terminating null byte is placed at end of buffer
  - *count*: number of bytes to read
    - (*buffer* must be at least this big)
    - (*size\_t* and *ssize\_t* are integer types)
- Returns:
  - $> 0$ : number of bytes read
    - May be  $< count$  (e.g., terminal *read()* gets only one line)
  - 0: end of file
  - -1: error



## *write()*: writing to a file

---

```
#include <unistd.h>
ssize_t write(int fd, const void *buffer, size_t count);
```

- Arguments:
  - *fd*: file descriptor
  - *buffer*: pointer to data to be written
  - *count*: number of bytes to write
- Returns:
  - Number of bytes written
    - May be less than *count* (e.g., device full, or insufficient space to write entire buffer to nonblocking socket)
  - $-1$  on error

## *close()*: closing a file

---

```
#include <unistd.h>
int close(fd);
```

- *fd*: file descriptor
- Returns:
  - 0: success
  - $-1$ : error
- Really should check for error!
  - Accidentally closing same FD twice
    - I.e., detect program logic error
  - Filesystem-specific errors
    - E.g., NFS commit failures may be reported only at *close()*
- **Note:** *close()* **always** releases FD, even on failure return
  - See *close(2)* man page

## Example: copy.c

---

```
$ ./copy old-file new-file
```

- A simple version of *cp(1)*

## Example: fileio/copy.c (snippet)

---

### Always remember to handle errors!

```
#define BUF_SIZE 1024
char buf[BUF_SIZE];

infd = open(argv[1], O_RDONLY);
if (infd == -1) errExit("open %s", argv[1]);

flags = O_CREAT | O_WRONLY | O_TRUNC;
mode = S_IRUSR | S_IWUSR | S_IRGRP; /* rw-r----- */
outfd = open(argv[2], flags, mode);
if (outfd == -1) errExit("open %s", argv[2]);

while ((nread = read(infd, buf, BUF_SIZE)) > 0)
    if (write(outfd, buf, nread) != nread)
        fatal("write() returned error or partial write occurred");
if (nread == -1) errExit("read");

if (close(infd) == -1) errExit("close");
if (close(outfd) == -1) errExit("close");
```

# Universality of I/O

---

- The fundamental I/O system calls work on almost all file types:

```
$ ls > mylist
$ ./copy mylist new          # Regular file

$ ./copy mylist /dev/tty    # Device

$ mkfifo f; cat f &         # FIFO
$ ./copy mylist f
```

## Notes for online practical sessions

---

- Small groups in **breakout rooms**
  - Write a note into Slack if you have a preferred group
- **We will go faster, if groups collaborate** on solving the exercise(s)
  - You can **share a screen** in your room
- I will circulate regularly between rooms to answer questions
- Zoom has an “**Ask for help**” button...
- **Keep an eye on the #general Slack channel**
  - Perhaps with further info about exercise;
  - Or a note that the exercise merges into a break
- When your room has finished, write a message in the Slack channel: “**\*\*\*\*\* Room X has finished \*\*\*\*\***”
  - Then I have an idea of how many people have finished

## Exercise notes

---

- For many exercises, there are templates for the solutions
  - Filenames: `ex.*.c`
  - Look for FIXMEs to see what pieces of code you must add
  - ⚠ You will need to edit the corresponding Makefile to add a new target for the executable
    - Look for the EXERCISE\_SOLNS\_EXE macro

```
-EXERCISE_FILES_EXE = # ex.prog_a ex.prog_b  
+EXERCISE_FILES_EXE = ex.prog_a # ex.prog_b
```

- Get a *make* tutorial now if you need one

## Exercise

---

- ① Using `open()`, `close()`, `read()`, and `write()`, implement the command `tee [-a] file` ([**template: fileio/ex.tee.c**]). This command writes a copy of its standard input to standard output and to `file`. If `file` does not exist, it should be created. If `file` already exists, it should be truncated to zero length (`O_TRUNC`). The program should support the `-a` option, which appends (`O_APPEND`) output to the file if it already exists, rather than truncating the file. Some hints:
  - Build `../libtspi.a` by doing *make* in source code root directory!
  - After first doing some simple command-line testing, test using the unit test in the Makefile: `make tee_test`.
  - Remember that you will need to add a target in the Makefile!
  - Standard input & output are automatically opened for a process.
  - Why does “`man open`” show the wrong manual page? It finds a page in the wrong section first. Try “`man 2 open`” instead.
  - `while inotifywait -q . ; do echo -e '\n\n'; make; done`
    - You may need to install the *inotify-tools* package
  - Command-line options can be parsed using `getopt(3)`.

# Outline

---

3	File I/O and Files	3-1
3.1	File I/O overview	3-3
3.2	<code>open()</code> , <code>read()</code> , <code>write()</code> , and <code>close()</code>	3-8
3.3	Retrieving file information: <code>stat()</code>	3-23

## Retrieving file information: `stat()`

---

```
#include <sys/stat.h>
int stat(const char *pathname, struct stat *statbuf);
int lstat(const char *pathname, struct stat *statbuf);
int fstat(int fd, struct stat *statbuf);
```

- Retrieve information about a file (“metadata”), mostly from inode
  - Information placed in *statbuf*
- `stat()`: retrieve info about **filename** identified by *pathname*
- `lstat()`: if *pathname* is a **symbolic link**, retrieve information about link, not file to which it refers
  - (`stat()` dereferences symbolic links)
- `fstat()`: retrieve info about file referred to by **descriptor** *fd*

[TLPI §15.1]

## The *stat* structure

```
struct stat {
    dev_t      st_dev;      /* ID of device containing file */
    ino_t      st_ino;      /* Inode number of file */
    mode_t     st_mode;     /* File type and permissions */
    nlink_t    st_nlink;    /* # of (hard) links to file */
    uid_t      st_uid;      /* User ID of file owner */
    gid_t      st_gid;      /* Group ID of file owner */
    dev_t      st_rdev;     /* ID for device special files */
    off_t      st_size;     /* File size (bytes) */
    blksize_t  st_blksize;  /* Optimal I/O block size (B) */
    blkcnt_t   st_blocks;   /* Storage allocated to file,
                           measured in 512B units */
    time_t     st_atime;    /* Time of last file access */
    time_t     st_mtime;    /* Time of last file modification */
    time_t     st_ctime;    /* Time of last change to inode */
};
```

- All types above are defined by POSIX (mostly integers)
- Full details on fields can be found in *inode(7)* and *stat(2)*
  - We'll look at details of a **subset** of these fields

## The *stat* structure

- *st\_dev*: ID of device containing filesystem where device resides
  - Consists of major ID (12 bits) + minor ID (20 bits)
  - *st\_dev* value is calculated by kernel (not stored as part of inode)
- *st\_ino*: inode number of file
  - Each file on a filesystem has a unique inode number
- *st\_nlink*: number of (hard) links to file
- *st\_size*: nominal file size (bytes) (*ls -l*)

## File timestamps

---

- File timestamps record time since Epoch (00:00:00, 1 Jan 1970, UTC):
  - *st\_atime*: time of last access of file data
  - *st\_mtime*: time of last modification of file data
  - *st\_ctime*: time of last change to inode
- Various system calls update timestamps as expected
  - TLPI Table 15-2
- In reality, timestamps are *timespec* structures, recording **seconds and nanoseconds**
  - E.g., *st\_atim.tv\_sec* and *st\_atim.tv\_nsec*
  - Not all FS types support nanosecond timestamps
    - XFS, ext4, and Btrfs do

[TLPI §15.2]

## File ownership

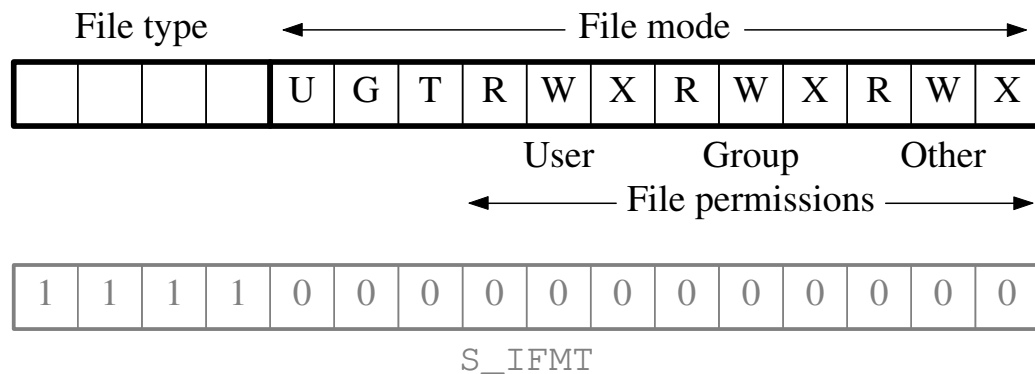
---

- *st\_uid* and *st\_gid* identify **ownership of file**
  - File UID + GID determine **permissions for file access**
- **UID of new file** == effective UID of creating process
- On most filesystems, **GID of new file** is either:
  - Effective **GID of creating process** (System V semantics)
  - **GID of parent directory** (BSD semantics)
    - Allows creation of subtrees that are always accessible to a particular group
- Choice is determined by whether **parent directory's set-GID bit** is enabled
  - `chmod g+s <dir>` (Propagates directory GID to new files)
  - This use of set-GID bit is a Linux extension

[TLPI §15.3]

## File type and mode

- `st_mode` returns two pieces of info:



- Left-most bits give **file type**
- Remaining bits are **file mode**
  - File permissions (9 bits) + set-UID/set-GID/sticky bits
- `S_IFMT` (defined in `<sys/stat.h>`) macro allows file type to be extracted

[TLPI §15.1]

## File type

- Extract using `statbuf.st_mode` & `S_IFMT`
- Check using predefined constants and macros:

```
if ((statbuf.st_mode & S_IFMT) == S_IFREG)
    printf("regular file\n");
```

- Common operation, so there are shorthand macros:

```
if (S_ISREG(statbuf.st_mode)) ...
```

Constant	Test macro	File type
S_IFREG	S_ISREG()	Regular file
S_IFDIR	S_ISDIR()	Directory
S_IFCHR	S_ISCHR()	Character device
S_IFBLK	S_ISBLK()	Block device
S_IFIFO	S_ISFIFO()	FIFO
S_IFSOCK	S_ISSOCK()	Socket
S_IFLNK	S_ISLNK()	Symbolic link



Notes

Notes

# Linux System Programming Essentials

# Processes

Michael Kerrisk, man7.org © 2021

mtk@man7.org

November 2021

## Outline

---

4	Processes	4-1
4.1	Process IDs	4-3
4.2	Process memory layout	4-6
4.3	Command-line arguments	4-9
4.4	The environment list	4-11
4.5	The /proc filesystem	4-16

# Outline

---

4	Processes	4-1
4.1	Process IDs	4-3
4.2	Process memory layout	4-6
4.3	Command-line arguments	4-9
4.4	The environment list	4-11
4.5	The /proc filesystem	4-16

## Process ID

---

```
#include <unistd.h>
pid_t getpid(void);
```

- **Process** == running instance of a program
  - Program + program loader (kernel) ⇒ process
- Every process has a process ID (PID)
  - *pid\_t*: positive integer that uniquely identifies process
  - *getpid()* returns callers's PID
  - Maximum PID is 32767 on Linux
    - Kernel then cycles, reusing PIDs, starting at low numbers
    - All PID slots used? ⇒ *fork()* fails with EAGAIN
    - Limit adjustable via /proc/sys/kernel/pid\_max (up to kernel's PID\_MAX\_LIMIT constant, typically 4\*1024\*1024)
    - (Process creation is also limited by RLIMIT\_NPROC process resource limit; see *getrlimit(2)*)

[TLPI §6.2]

# Parent process ID

---

```
#include <unistd.h>
pid_t getppid(void);
```

- Every process has a parent
  - Typically, process that created this process using *fork()*
  - Parent process is informed when its child terminates
- All processes on system thus form a tree
  - At root is *init*, PID 1, the ancestor of all processes
  - “Orphaned” processes are “adopted” by *init*
- *getppid()* returns PID of caller’s parent process (PPID)

[TLPI §6.2]

## Outline

---

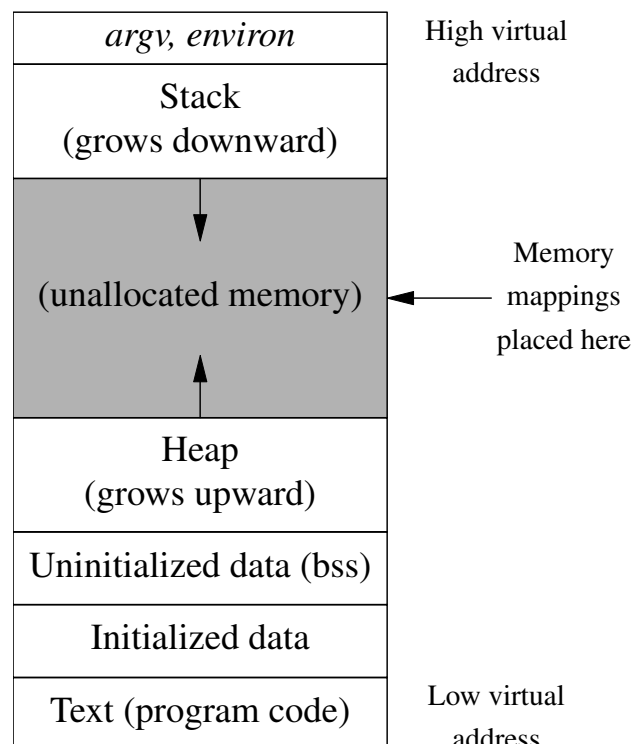
4	Processes	4-1
4.1	Process IDs	4-3
4.2	Process memory layout	4-6
4.3	Command-line arguments	4-9
4.4	The environment list	4-11
4.5	The /proc filesystem	4-16

# Process memory layout

Virtual memory of a process is divided into **segments**:

- **Text**: machine-language instructions
  - Marked read-only to prevent self-modification
  - Multiple processes can share same code in memory
- **Initialized data**: global and static variables that are explicitly initialized
  - Values read from program file when process is created
- **Uninitialized data**: global and static variables that are not explicitly initialized
  - Initialized to zero when process is created
- **Stack**: storage for function local variables and call linkage info (saved SP and PC registers)
- **Heap**: an area from which memory can be dynamically allocated and deallocated
  - *malloc()* and *free()*

## Process memory layout (simplified)



[TLPI §6.3]

# Outline

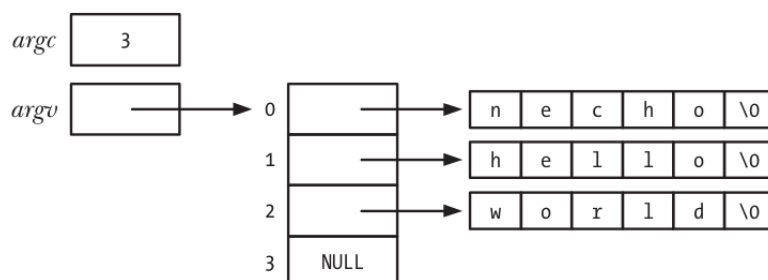
---

4	Processes	4-1
4.1	Process IDs	4-3
4.2	Process memory layout	4-6
4.3	Command-line arguments	4-9
4.4	The environment list	4-11
4.5	The /proc filesystem	4-16

## Command-line arguments

---

- Command-line arguments of a program provided as first two arguments of *main()*
  - Conventionally named *argc* and *argv*
- *int argc*: number of arguments
- *char \*argv[]*: array of pointers to arguments (strings)
  - *argv[0]* == name used to invoke program
  - *argv[argc]* == *NULL*
- E.g., for the command, *necho hello world*:



[TLPI §6.6]

# Outline

---

4	Processes	4-1
4.1	Process IDs	4-3
4.2	Process memory layout	4-6
4.3	Command-line arguments	4-9
4.4	The environment list	4-11
4.5	The /proc filesystem	4-16

## Environment list (*environ*)

---

Each process has a list of **environment variables**

- Strings of form *name=value*
- New process inherits copy of parent's environment
  - Simple (one-way) interprocess communication
- Commonly used to control behavior of programs
- Examples:
  - HOME: user's home directory (initialized at login)
  - PATH: list of directories to search for executable programs
  - EDITOR: user's preferred editor

[TLPI §6.7]

## Environment list (*environ*)

- Can create environment variables within shell:

```
$ MANWIDTH=72
$ export MANWIDTH
$ man getpid
```

- All processes created by shell will inherit definition
- Creating an environment variable for a single command (does not modify shell's environment):

```
$ MANWIDTH=72 man getpid
```

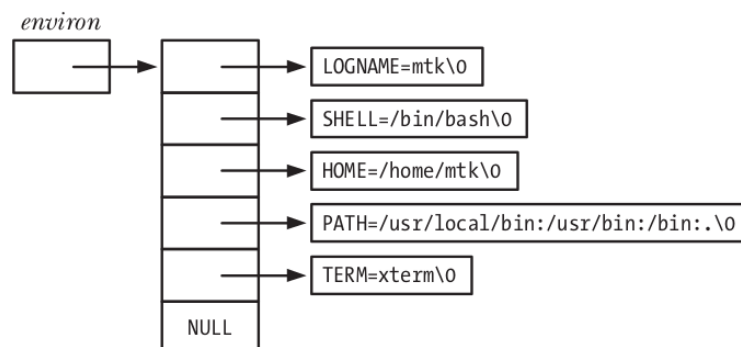
- To list all environment variables, use *env(1)* or *printenv(1)*

## Accessing the environment from a program

- Environment list can be accessed via a global variable:

```
extern char **environ;
```

- NULL-terminated array of pointers to strings:



- Displaying environment:

```
for (char **ep = environ; *ep != NULL; ep++)
    puts(*ep);
```



## Environment variable APIs

---

- Fetching value of an EV: `value = getenv("NAME");`
- Creating/modifying an EV:
  - `putenv("NAME=value");`
  - `setenv("NAME", "value", overwrite);`
- Removing an EV: `unsetenv("NAME");`
- `/proc/PID/envIRON` can be used (with suitable permissions) to view environment of another process
- See man pages and TLPI §6.7

## Outline

---

4	Processes	4-1
4.1	Process IDs	4-3
4.2	Process memory layout	4-6
4.3	Command-line arguments	4-9
4.4	The environment list	4-11
4.5	The <code>/proc</code> filesystem	4-16

## The /proc filesystem

---

- Pseudofilesystem that exposes kernel information via filesystem metaphor
  - Structured as a set of subdirectories and files
  - *proc(5)* man page
- Files don't really exist
  - Created on-the-fly when pathnames under /proc are accessed
- Many files read-only
- Some files are writable  $\Rightarrow$  can update kernel settings

## The /proc filesystem: examples

---

- /proc/cmdline: command line used to start kernel
- /proc/cpuinfo: info about CPUs on the system
- /proc/meminfo: info about memory and memory usage
- /proc/modules: info about loaded kernel modules
- /proc/sys/fs/: files and subdirectories with filesystem-related info
- /proc/sys/kernel/: files and subdirectories with various readable/settable kernel parameters
- /proc/sys/net/: files and subdirectories with various readable/settable networking parameters

`/proc/PID/` directories

- One `/proc/PID/` subdirectory for each running process
- Subdirectories and files exposing info about process with corresponding PID
- Some files publicly readable, some readable only by process owner; a few files writable
- Examples
  - `cmdline`: command line used to start program
  - `cwd`: current working directory
  - `environ`: environment of process
  - `fd`: directory with info about open file descriptors
  - `limits`: resource limits
  - `maps`: mappings in virtual address space
  - `status`: (lots of) info about process

## Notes

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

# Linux System Programming Essentials

# Signals

Michael Kerrisk, man7.org © 2021

mtk@man7.org

November 2021

## Outline

---

5	Signals	5-1
5.1	Overview of signals	5-3
5.2	Signal dispositions	5-8
5.3	Useful signal-related functions	5-16
5.4	Signal handlers	5-20
5.5	Signal sets, the signal mask, and pending signals	5-25
5.6	Designing signal handlers	5-37

# Outline

---

5	Signals	5-1
5.1	Overview of signals	5-3
5.2	Signal dispositions	5-8
5.3	Useful signal-related functions	5-16
5.4	Signal handlers	5-20
5.5	Signal sets, the signal mask, and pending signals	5-25
5.6	Designing signal handlers	5-37

## Signals are a notification mechanism

---

- Signal == notification to a process that an event occurred
  - “Software interrupts”
  - **asynchronous**: receiver (generally) can’t predict when a signal will occur

## Signal types

---

- 64 signals (on Linux)
- Each signal has a unique integer value
  - Numbered starting at 1 ⚠
- Defined symbolically in `<signal.h>`:
  - Names of form `SIGxxx`
  - e.g., signal 2 is `SIGINT` (“terminal interrupt”)
- Two broad categories of signals:
  - “Standard” signals (1 to 31)
    - Mostly for kernel-defined purposes
  - Realtime signals (32 to 64)
    - Exist for user-defined purposes

[TLPI §20.1]

## Signal generation

---

- Signals can be sent by:
  - The kernel (the common case)
  - Another process (with suitable permissions)
    - `kill(pid, sig)` and related APIs
- Kernel generates signals for various events, e.g.:
  - Attempt to access a nonexistent memory address (`SIGSEGV`)
  - Terminal *interrupt* character (Control-C) was typed (`SIGINT`)
  - Child process terminated (`SIGCHLD`)
  - Process CPU time limit exceeded (`SIGXCPU`)

[TLPI §20.1]

# Terminology

---

Some terminology:

- A signal is **generated** when an event occurs
- Later, a signal is **delivered** to the process, which then takes some action in response
- Between generation and delivery, a signal is **pending**
- We can **block** (delay) delivery of specific signals by adding them to process's **signal mask**
  - **Signal mask == set of signals whose delivery is blocked**
  - Pending signal is delivered only after it is unblocked

[TLPI §20.1]

## Outline

---

5	Signals	5-1
5.1	Overview of signals	5-3
5.2	Signal dispositions	5-8
5.3	Useful signal-related functions	5-16
5.4	Signal handlers	5-20
5.5	Signal sets, the signal mask, and pending signals	5-25
5.6	Designing signal handlers	5-37

# Signal default actions

- When a signal is delivered, a process takes one of these default actions:
  - **Ignore**: signal is discarded by kernel, has no effect on process
  - **Terminate**: process is terminated (“killed”)
  - **Core dump**: process produces a core dump and is terminated
    - Core dump file can be used to examine state of program inside a debugger
    - See also *core(5)* man page
  - **Stop**: execution of process is suspended
  - **Continue**: execution of a stopped process is resumed
- Default action for each signal is signal-specific

[TLPI §20.2]

# Standard signals and their default actions

Name	Description	Default
SIGABRT	Abort process	Core
SIGALRM	Real-time timer expiration	Term
SIGBUS	Memory access error	Core
SIGCHLD	Child stopped or terminated	Ignore
SIGCONT	Continue if stopped	Cont
SIGFPE	Arithmetic exception	Core
SIGHUP	Hangup	Term
SIGILL	Illegal Instruction	Core
SIGINT	Interrupt from keyboard	Term
SIGIO	I/O Possible	Term
SIGKILL	Sure kill	Term
SIGPIPE	Broken pipe	Term
SIGPROF	Profiling timer expired	Term
SIGPWR	Power about to fail	Term
SIGQUIT	Terminal quit	Core
SIGSEGV	Invalid memory reference	Core
SIGSTKFLT	Stack fault on coprocessor	Term
SIGSTOP	Sure stop	Stop
SIGSYS	Invalid system call	Core
SIGTERM	Terminate process	Term
SIGTRAP	Trace/breakpoint trap	Core
SIGTSTP	Terminal stop	Stop
SIGTTIN	Terminal input from background	Stop
SIGTTOU	Terminal output from background	Stop
SIGURG	Urgent data on socket	Ignore
SIGUSR1	User-defined signal 1	Term
SIGUSR2	User-defined signal 2	Term
SIGVTALRM	Virtual timer expired	Term
SIGWINCH	Terminal window size changed	Ignore
SIGXCPU	CPU time limit exceeded	Core
SIGXFSZ	File size limit exceeded	Core

- Signal default actions are:
  - Term: terminate the process
  - Core: produce core dump and terminate the process
  - Ignore: ignore the signal
  - Stop: stop (suspend) the process
  - Cont: resume process (if stopped)
- SIGKILL and SIGSTOP can't be caught, blocked, or ignored
- TLPI §20.2



## Stop and continue signals

---

- Certain signals **stop** a process, freezing its execution
- Examples:
  - SIGTSTP: “terminal stop” signal, generated by typing Control-Z
  - SIGSTOP: “sure stop” signal
- SIGCONT causes a stopped process to resume execution
  - SIGCONT is ignored if process is not stopped
- Most common use of these signals is in **shell job control**

## Changing a signal's disposition

---

- Instead of default, we can change a signal's disposition to:
  - **Ignore** the signal
  - **Handle (“catch”) the signal**: execute a user-defined function upon delivery of the signal
  - Revert to the **default action**
    - Useful if we earlier changed disposition
- Can't change disposition to *terminate* or *core dump*
  - But, a signal handler can emulate these behaviors
- Can't change disposition of SIGKILL or SIGSTOP (EINVAL)
  - So, they always kill or stop a process

## Changing a signal's disposition: *sigaction()*

```
#include <signal.h>
int sigaction(int sig,
               const struct sigaction *act,
               struct sigaction *oldact);
```

*sigaction()* changes (and/or retrieves) disposition of signal *sig*

- *sigaction* structure describes a signal's disposition
- *act* points to structure specifying new disposition for *sig*
  - Can be NULL for no change
- *oldact* returns previous disposition for *sig*
  - Can be NULL if we don't care
- *sigaction(sig, NULL, &oldact)* returns current disposition, without changing it

[TLPI §20.13]

## *sigaction* structure

```
struct sigaction {
    void (*sa_handler)(int);
    sigset_t sa_mask;
    int sa_flags;
    void (*sa_restorer)(void);
};
```

- *sa\_handler* specifies disposition of signal:
  - Address of a signal handler function
  - SIG\_IGN: ignore signal
  - SIG\_DFL: revert to default disposition
- *sa\_mask*: signals to block while handler is executing
  - Field is initialized using macros described in *sigsetops(3)*
- *sa\_flags*: bit mask of flags affecting invocation of handler
- *sa\_restorer*: not for application use
  - Used internally to implement “signal trampoline”

## Ignoring a signal (signals/ignore\_signal.c)

---

```
int ignoreSignal(int sig)
{
    struct sigaction sa;

    sa.sa_handler = SIG_IGN;
    sa.sa_flags = 0;
    sigemptyset(&sa.sa_mask);
    return sigaction(sig, &sa, NULL);
}
```

- A “library function” that ignores specified signal
- Other fields only significant when establishing a signal handler, but must be properly initialized here

## Outline

---

5	Signals	5-1
5.1	Overview of signals	5-3
5.2	Signal dispositions	5-8
5.3	Useful signal-related functions	5-16
5.4	Signal handlers	5-20
5.5	Signal sets, the signal mask, and pending signals	5-25
5.6	Designing signal handlers	5-37

# Displaying signal descriptions

```
#define _GNU_SOURCE
#include <string.h>
char *strsignal(int sig);
```

- Returns string describing signal *sig*
- NSIG constant is 1 greater than maximum signal number
  - Define `_GNU_SOURCE` to get definition from `<signal.h>`

[TLPI §20.8]

## Example: `signals/t_strsignal.c`

```
int main(int argc, char *argv[]) {
    for (int sig = 1; sig < NSIG; sig++)
        printf("%2d: %s\n", sig, strsignal(sig));

    exit(EXIT_SUCCESS);
}
```

```
$ ./t_strsignal
1: Hangup
2: Interrupt
3: Quit
4: Illegal instruction
5: Trace/breakpoint trap
6: Aborted
7: Bus error
8: Floating point exception
9: Killed
10: User defined signal 1
11: Segmentation fault
12: User defined signal 2
13: Broken pipe
...
```

## Waiting for a signal: *pause()*

---

```
#include <unistd.h>
int pause(void);
```

- Blocks execution of caller until a signal is caught
- Always returns `-1` with *errno* set to `EINTR`
  - (Standard return for blocking system call that is interrupted by a signal handler)

[TLPI §20.14]

## Outline

---

5	Signals	5-1
5.1	Overview of signals	5-3
5.2	Signal dispositions	5-8
5.3	Useful signal-related functions	5-16
5.4	Signal handlers	5-20
5.5	Signal sets, the signal mask, and pending signals	5-25
5.6	Designing signal handlers	5-37

# Signal handlers

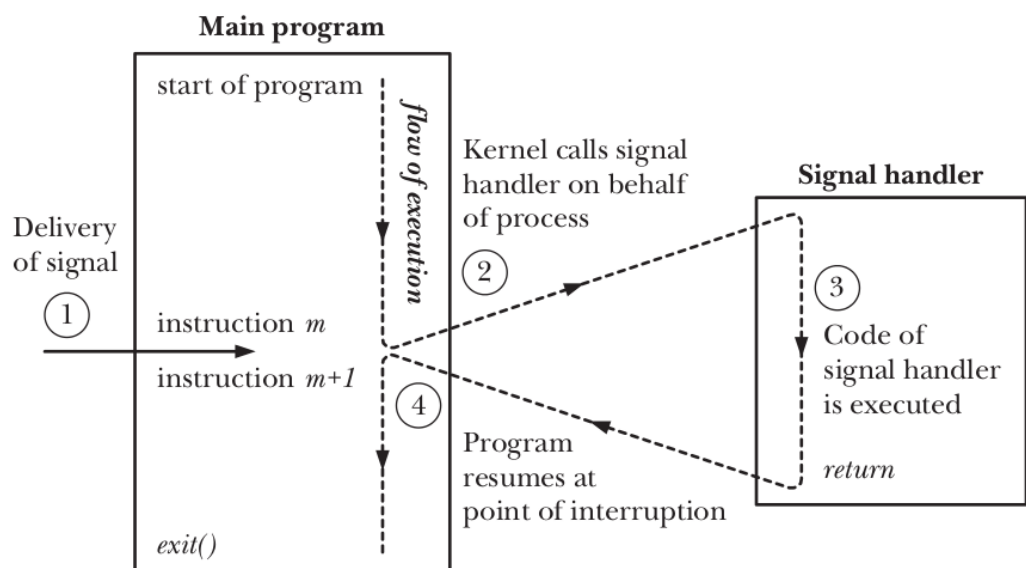
- Programmer-defined function
- Called with one integer argument: number of signal
  - $\Rightarrow$  handler installed for multiple signals can differentiate...
- Returns void

```
void
myHandler(int sig)
{
    /* Actions to be performed when signal
       is delivered */
}
```

[TLPI §20.4]

## Signal handler invocation

- Automatically invoked by kernel when signal is delivered:
  - Can interrupt main program flow at any time
  - On return, execution continues at point of interruption



## Example: signals/ouch\_sigaction.c

Print "Ouch!" when Control-C is typed at keyboard

```
static void sigHandler(int sig) {
    printf("Ouch!\n");          /* UNSAFE */
}

int main(int argc, char *argv[]) {
    struct sigaction sa;
    sa.sa_flags = 0;           /* No flags */
    sa.sa_handler = sigHandler; /* Handler function */
    /* Don't block additional signals
       during invocation of handler */
    sigemptyset(&sa.sa_mask);

    if (sigaction(SIGINT, &sa, NULL) == -1)
        errExit("sigaction");

    for (;;)
        pause();               /* Wait for a signal */
}
```

## Exercise

- While a signal handler is executing, the signal that caused it to be invoked is (by default) temporarily added to the signal mask, so that it is blocked from further delivery until the signal handler returns. Consequently, execution of a signal handler can't be interrupted by a further execution of the same handler. To demonstrate that this is so, modify the signal handler in the signals/ouch\_sigaction.c program to include the following after the existing *printf()* statement:

```
sleep(5);
printf("Bye\n");
```

Build and run the program, type control-C once, and then while the signal handler is executing, type control-C three more times. What happens? In total, how many times is the signal handler called?

# Outline

---

5	Signals	5-1
5.1	Overview of signals	5-3
5.2	Signal dispositions	5-8
5.3	Useful signal-related functions	5-16
5.4	Signal handlers	5-20
5.5	Signal sets, the signal mask, and pending signals	5-25
5.6	Designing signal handlers	5-37

## Signal sets

---

- Various signal-related APIs work with **signal sets**
- Signal set == data structure that represents multiple signals
- Data type: *sigset\_t*
  - Typically a bit mask, but not necessarily



# Manipulating signal sets

```
#include <signal.h>
int sigemptyset(sigset_t *set);
int sigfillset(sigset_t *set);
int sigaddset(sigset_t *set, int sig);
int sigdelset(sigset_t *set, int sig);
int sigismember(const sigset_t *set, int sig);
```

- *sigemptyset()* initializes *set* to contain no signals
- *sigfillset()* initializes *set* to contain all signals
  - We **must** initialize *set* using *sigemptyset()* or *sigfillset()* before employing macros below
  - Using *memset()* to zero a signal set is *not* correct
- *sigaddset()* adds *sig* to *set*
- *sigdelset()* removes *sig* from *set*
- *sigismember()* returns 1 if *sig* is in *set*, 0 if it is not, or -1 on error (e.g., *sig* is invalid)

## Blocking signals (the signal mask)

- Each process has a **signal mask**—a set of signals whose delivery is currently blocked
  - (In truth: each **thread** has a signal mask...)
- If a blocked signal is generated, it remains pending until removed from signal mask
- The signal mask can be changed in various ways:
  - While handler is invoked, the **signal that triggered the handler** is (temporarily) added to signal mask
  - While handler is invoked, any signals specified in *sa\_mask* are (temporarily) added to signal mask
  - Explicitly, using *sigprocmask()*
- Attempts to block SIGKILL/SIGSTOP are silently ignored

## *sigprocmask()*

---

```
#include <signal.h>
int sigprocmask(int how, const sigset_t *set,
                sigset_t *oldset);
```

- Adds signals to, or removes signals from, caller's signal mask
  - (Typical use: prevent interruption by signal handler while updating a shared data structure)
- *how* specifies change to signal mask:
  - SIG\_BLOCK: **add** signals in *set* to signal mask
  - SIG\_UNBLOCK: **remove** signals in *set* from signal mask
  - SIG\_SETMASK: **assign** *set* to signal mask

[TLPI §20.10]

## *sigprocmask()*

---

```
#include <signal.h>
int sigprocmask(int how, const sigset_t *set,
                sigset_t *oldset);
```

- *oldset* returns previous signal mask
  - Can be NULL if we don't care
- `sigprocmask(how, NULL, &oldset)` retrieves current mask without changing it
  - *how* is ignored

[TLPI §20.10]

## Example: temporarily blocking a signal

- The following code snippet shows how to temporarily block a signal (SIGINT) while executing a block of code

```
sigset_t blocking, prev;

sigemptyset(&blocking);
sigaddset(&blocking, SIGINT);
sigprocmask(SIG_BLOCK, &blocking, &prev);

/* ... Code to execute with SIGINT blocked ... */

sigprocmask(SIG_SETMASK, &prev, NULL);
```

## Pending signals

```
#include <signal.h>
int sigpending(sigset_t *set);
```

- Between generation and delivery, a signal is **pending**
  - Pending state is normally unobservable unless signal is explicitly blocked
- *sigpending()* returns (in *set*) the set of signals currently pending for caller
  - We do **not** need to initialize *set* before calling *sigpending()*
- Can examine *set* using *sigismember()*:

```
sigset_t pending;
sigpending(&pending);
if (sigismember(&pending, SIGINT))
    printf("SIGINT (%s) is pending\n",
           strsignal(SIGINT));
```

## Signals are not queued

---

- The set of pending (standard) signals is a mask
- $\Rightarrow$  If same signal is generated multiple times while blocked, it will be delivered just once

## Exercises

---

The goal of this exercise is experiment with signal handlers and the use of the signal mask to block delivery of signals. A template for the complete exercise is provided (**[template: signals/ex.pending\_sig\_expt.c]**)

**Hint:** don't confuse the *sa\_mask* field that is passed to *sigaction()*, which specifies additional signals that should be temporarily blocked while a signal handler is executing, with the use of *sigprocmask()*, which allows a process to directly modify its signal mask.

① Write a program that:

- Blocks all signals except SIGINT (*sigprocmask()*, slides 5-30 + 5-31).
- Uses *sigaction()* (slides 5-13 + 5-14) to establish a SIGINT handler that does nothing but return.
- Calls *pause()* to wait for a signal.

[Exercises continue on following slides]

## Exercises

---

- After *pause()* returns, determines the set of pending signals for the process (use *sigpending()*, slide 5-32), tests which signals are in that set (use *sigismember()*, iterating through all signals in the range  $1 \leq s < \text{NSIG}$ ; see slide 5-18), and prints their descriptions (*strsignal()*).

Run the program and send it various signals (other than SIGINT and signals that are ignored by default), using either the *kill* command (*kill -<sig> <pid>*) from another terminal, or by typing signal-generating keys from the terminal where you run the program (Control-Z for SIGTSTP, Control-\ for SIGQUIT). Then type Control-C to generate SIGINT and inspect the list of pending signals displayed by the program.

- ② What happens if you send SIGKILL to the preceding program? Why?

[Exercises continue on following slide]

## Exercises

---

- ③ Extend the program created in the preceding exercise so that:
- Just after installing the handler for SIGINT, the program installs an additional handler for SIGQUIT (generated when the Control-\ key is pressed). The handler should print a message "SIGQUIT received", and return.
  - After displaying the list of pending signals, the program unblocks SIGQUIT and calls *pause()* once more. ( ⚠ Which *how* value should be given to *sigprocmask()*?)

While the program is blocking signals (i.e., before typing Control-C), try typing Control-\ multiple times. After Control-C is typed, how many times does the SIGQUIT handler display its message? Why?

## Outline

---

5	Signals	5-1
5.1	Overview of signals	5-3
5.2	Signal dispositions	5-8
5.3	Useful signal-related functions	5-16
5.4	Signal handlers	5-20
5.5	Signal sets, the signal mask, and pending signals	5-25
5.6	Designing signal handlers	5-37

## Keep it simple

---

- Signal handlers can, in theory, do anything
- But, complex signal handlers can easily have subtle bugs (e.g., race conditions)
  - E.g., if main program and signal handler access same global variables
- ⇒ Avoid using signals if you can
  - ⚠ Don't introduce them as a means of IPC
  - ⚠ Don't use as part of a library design
    - (That would imply a contract with main program about which signals library is allowed to use)
- But, in some cases, we must deal with signals sent by kernel
  - ⇒ Design the handlers to be as simple as possible

## Keep it simple

---

- Some simple signal-handler designs:
  - Set a global flag and return
    - Main program periodically checks (and clears) flag, and takes appropriate action
  - Signal handler does some clean-up and terminates process
    - (TLPI §21.2)
  - Signal handler performs a nonlocal goto to unwind stack
    - *sigsetjmp()* and *siglongjmp()* (TLPI §21.2.1)
    - E.g., some shells do this when handling signals

## Signals are not queued

---

- Signals are not queued
- A blocked signal is marked just once as pending, even if generated multiple times
- ⇒ **One signal may correspond to multiple “events”**
  - Programs that handle signals must be designed to allow for this
- Example:
  - SIGCHLD is generated for parent when child terminates
  - While SIGCHLD handler executes, SIGCHLD is blocked
  - Suppose **two** more children terminate while handler executes
  - Only one SIGCHLD signal will be queued
  - Solution: SIGCHLD handler should loop, checking if multiple children have terminated

Linux System Programming Essentials

# Process Lifecycle

Michael Kerrisk, man7.org © 2021

mtk@man7.org

November 2021

## Outline

---

6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: <code>fork()</code>	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The <code>SIGCHLD</code> signal	6-38
6.7	Executing programs: <code>execve()</code>	6-42



# Outline

---

6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: <code>fork()</code>	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The <code>SIGCHLD</code> signal	6-38
6.7	Executing programs: <code>execve()</code>	6-42

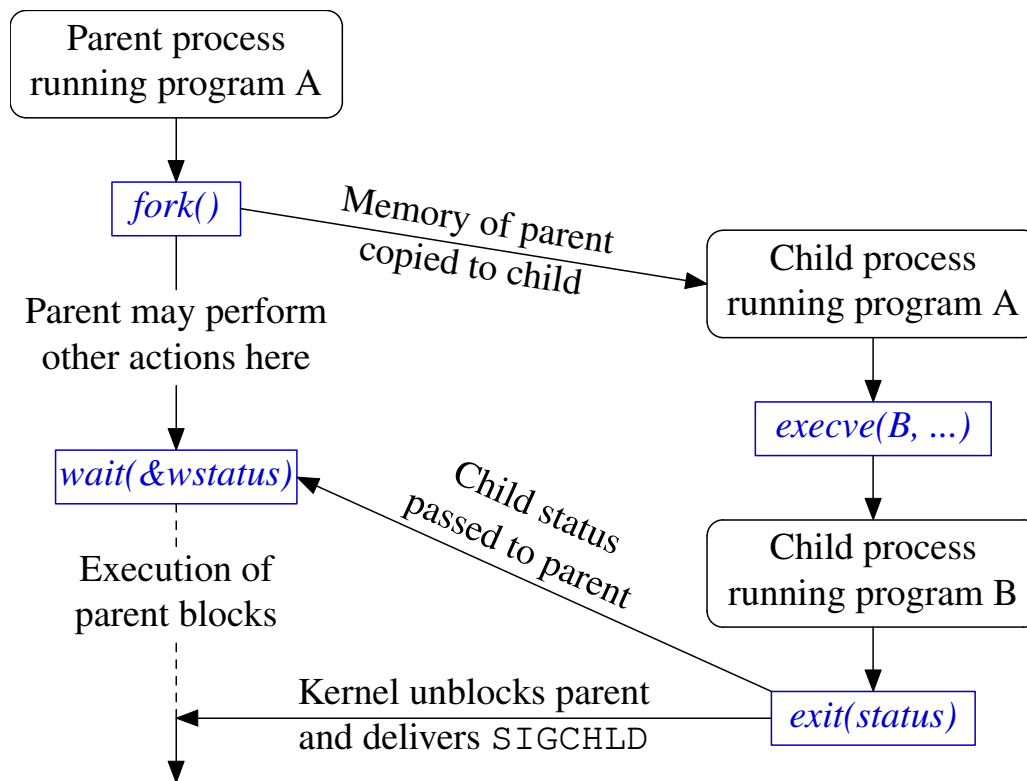
## Creating processes and executing programs

---

Four key system calls (and their variants):

- *fork()*: create a new (“child”) process
- *exit()*: terminate calling process
- *wait()*: wait for a child process to terminate
- *execve()*: execute a new program in calling process

## Using *fork()*, *execve()*, *wait()*, and *exit()* together



## Outline

6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: <code>fork()</code>	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The <code>SIGCHLD</code> signal	6-38
6.7	Executing programs: <code>execve()</code>	6-42

## Creating a new process: *fork()*

```
#include <unistd.h>
pid_t fork(void);
```

*fork()* creates a new process (“**the child**”):

- Child is a **near exact duplicate of caller** (“the parent”)
- Notionally, memory of parent is duplicated to create child
  - In practice, copy-on-write duplication is used
    - $\Rightarrow$  Only page tables must be duplicated at time of *fork()*
- Two processes share same (read-only) text segment
- Two processes have separate copies of stack, data, and heap segments
  - $\Rightarrow$  Each process can modify variables without affecting other process

[TLPI §24.2]

## Return value from *fork()*

```
#include <unistd.h>
pid_t fork(void);
```

- **Both** processes continue execution by returning from *fork()*
- *fork()* returns different values in parent and child:
  - Parent:
    - On success: PID of new child (allows parent to track child)
    - On failure: -1
  - Child: returns 0
    - Child can obtain its own PID using *getpid()*
    - Child can obtain PID of parent using *getppid()*

## Using *fork()*

---

```
pid_t pid = fork();
if (pid == -1) {
    /* Handle error */ ;
} else if (pid == 0) {
    /* Code executed by child */
} else {
    /* Code executed by parent */
}
```

## A Linux-specific alternative: *clone()*

---

- *clone()/clone3()* is another way of creating a process
- Much more flexibility than *fork()* (multiple arguments)
- Features include:
  - Parent and child may share various attributes (threads!)
    - Process ID
    - File descriptors
    - Virtual address space
    - Signal dispositions
  - Create new namespaces
  - Can obtain PID file descriptor that refers to child
    - Can wait/signal via PID FD
- Used to implement *pthread\_create()* (and, in glibc, *fork()*!)

## Exercise

---

- 1 Write a program that uses *fork()* to create a child process ([**template:** `procexec/ex.fork_var_test.c`]). After the *fork()* call, both the parent and child should display their PIDs (*getpid()*). Include code to demonstrate that the child process created by *fork()* can modify its copy of a local variable in *main()* without affecting the value in the parent's copy of the variable.

Note: you may find it useful to use the *sleep(num-secs)* library function to delay execution of the parent for a few seconds, to ensure that the child has a chance to execute before the parent inspects its copy of the variable.

## Exercise

---

- 2 Processes have many attributes. When a new process is created using *fork()*, which of those attributes are inherited by the child and which are not (e.g., are reset to some default)? Here, we explore whether two process attribute—signal dispositions and alarm timers—are inherited by a child process.

Write a program ([**template:** `procexec/ex.inherit_alarm.c`]) that performs the following steps in order to determine if a child process inherits signal dispositions and alarm timers from the parent:

- Establishes a SIGALRM handler that prints the process's PID.
- Starts an alarm timer that expires after two seconds. Do this using the call *alarm(2)*. When the timer expires, it will notify by sending a SIGALRM signal to the process.
- Creates a child process using *fork()*.
- After the *fork()*, the child fetches the disposition of the SIGALRM signal (*sigaction()*) and tests whether the *sa\_handler* field in the returned structure is the address of the signal handler
- Both processes then loop 5 times, sleeping for half a second (use *usleep()*) and displaying the process PID. Which of the processes receives a SIGALRM signal?

# Outline

---

6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: <code>fork()</code>	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The <code>SIGCHLD</code> signal	6-38
6.7	Executing programs: <code>execve()</code>	6-42

## Terminating a process

---

A process can terminate itself using two APIs:

- `_exit(2)` (system call)
- `exit(3)` (library function)

## Terminating a process with `_exit(2)`

```
#include <unistd.h>
void _exit(int status);
```

`_exit()` terminates the calling process

- AKA **normal termination**
  - **abnormal termination** == killed by a signal
- (In truth: on Linux, `_exit()` is a wrapper for Linux-specific `exit_group(2)`, which terminates all threads in a process)

## Process exit status

```
#include <unistd.h>
void _exit(int status);
```

- Least significant 8 bits of *status* define **exit status**
  - Remaining bits ignored
  - 0 == success
  - nonzero == failure
- POSIX specifies two constants:

```
#define EXIT_SUCCESS 0
#define EXIT_FAILURE 1
```

## Terminating a process with *exit(3)*

- Most programs employ *exit(3)*, rather than *\_exit(2)*

```
#include <stdlib.h>
void exit(int status);
```

- The *exit(3)* library function:
  - Calls exit handlers registered by process
    - Exit handler == callback function automatically called at normal process termination
    - *atexit(3)*, *on\_exit(3)*
  - Flushes *stdio* buffers
    - i.e., *\_exit()* does **not** flush *stdio* buffers
  - Calls: *\_exit(status)*
- return *n* inside *main()* is equivalent to *exit(n)*

## Process teardown

As part of process termination (normal or abnormal), various cleanups are performed:

- All open **file descriptors** are closed
  - Associated **file locks** are released
- Open **POSIX IPC objects** are closed (message queues, semaphores, shared memory)
- **Memory mappings** are unmapped
- **Memory locks** are removed
- **System V shared memory segments** are detached
- And more...



## Outline

---

6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: <code>fork()</code>	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The <code>SIGCHLD</code> signal	6-38
6.7	Executing programs: <code>execve()</code>	6-42

## Overview

---

- Parent processes can use the “wait” family of system calls to monitor state change events in child processes:
  - Termination
  - Stop (because of a signal)
  - Continue (after `SIGCONT` signal)
- Parent can obtain various info about state changes:
  - Exit status of process
  - What signal stopped or killed process
  - Whether process produced a core dump before terminating
- For historical reasons, there are multiple “wait” functions

## Waiting for children with *waitpid()*

```
#include <sys/wait.h>
pid_t waitpid(pid_t pid, int *wstatus, int options);
```

- *waitpid()* waits for a child process to change state
  - No child has changed state  $\Rightarrow$  call blocks
  - Child has already changed state  $\Rightarrow$  call returns immediately
- State change is reported in *wstatus* (if non-NULL)
  - (details later...)
- Return value:
  - On success: PID of child whose status is being reported
  - On error,  $-1$ 
    - No more children?  $\Rightarrow$  ECHILD

[TLPI §26.1.2]

## Waiting for children with *waitpid()*

```
#include <sys/wait.h>
pid_t waitpid(pid_t pid, int *wstatus, int options);
```

*pid* specifies which child(ren) to wait for:

- *pid* ==  $-1$ : **any** child of caller
- *pid* >  $0$ : child whose **PID** equals *pid*
- *pid* ==  $0$ : any child in **same process group** as caller
- *pid* <  $-1$ : any child in **process group whose ID equals *abs(pid)***
  - See *credentials(7)* and *setpgid(2)* for info on process groups

## Waiting for children with *waitpid()*

```
#include <sys/wait.h>
pid_t waitpid(pid_t pid, int *wstatus, int options);
```

- By default, *waitpid()* reports only **terminated** children
- The *options* bit mask can specify additional state changes to report:
  - WUNTRACED: report **stopped** children
  - WCONTINUED: report stopped children that have **continued**
- Specifying WNOHANG in *options* causes **nonblocking** wait
  - If no children have changed state, *waitpid()* returns immediately, with return value of 0

## *waitpid()* example

Wait for all children to terminate, and report their PIDs:

```
for (;;) {
    childPid = waitpid(-1, NULL, 0);
    if (childPid == -1) {
        if (errno == ECHILD) {
            printf("No more children!\n");
            break;
        } else {          /* Unexpected error */
            errExit("waitpid");
        }
    }

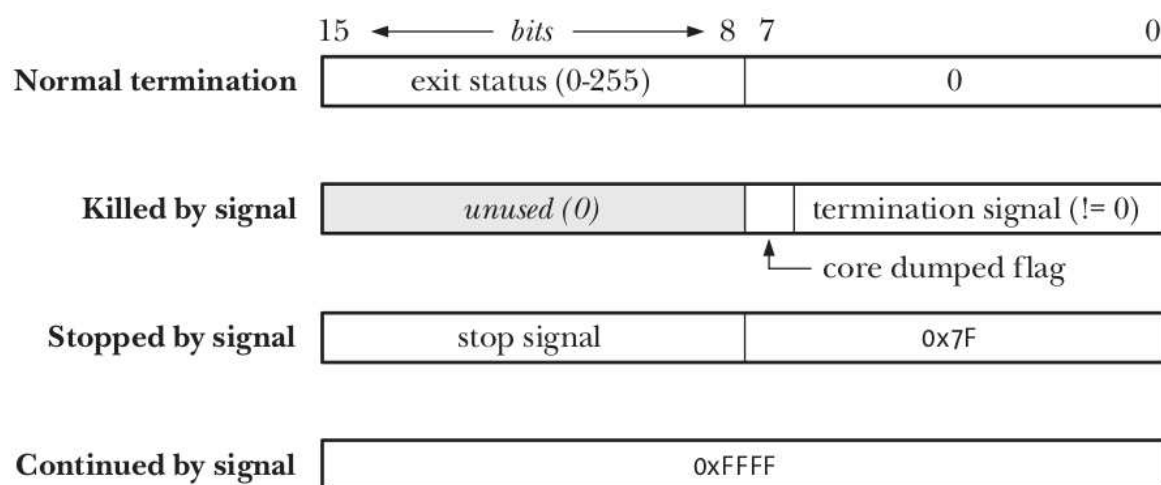
    printf("waitpid() returned PID %ld\n",
          (long) childPid);
}
```

## The wait status value

- *wstatus* returned by *waitpid()* distinguishes 4 types of event:
  - Child **terminated via `_exit()`**, specifying an *exit status*
  - Child was **killed by a signal**
  - Child was **stopped by a signal**
  - Child was **continued by a signal**
- The term **wait status** encompasses all four cases

## The wait status value

16 lowest bits of *wstatus* returned by *waitpid()* encode status in such a way that the 4 cases can be distinguished:



(Encoding is an implementation detail we don't really need to care about)

## Dissecting the wait status

- `<sys/wait.h>` defines macros for dissecting a wait status
- Only one of the headline macros in this list will return true:
  - ① `WIFEXITED(wstatus)`: true if child exited normally
    - `WEXITSTATUS(wstatus)` returns exit status of child
  - ② `WIFSIGNALED(wstatus)`: true if child was killed by signal
    - `WTERMSIG(wstatus)` returns number of killing signal
    - `WCOREDUMP(wstatus)` returns true if child dumped core
  - ③ `WIFSTOPPED(wstatus)`: true if child was stopped by signal
    - `WSTOPSIG(wstatus)` returns number of stopping signal
  - ④ `WIFCONTINUED(wstatus)`: true if child was resumed by `SIGCONT`
- The subordinate macros may be used only if the corresponding headline macro tests true

## Example: `procexec/print_wait_status.c`

Display wait status value in human-readable form

```
1 void printWaitStatus(const char *msg, int status) {
2     if (msg != NULL)
3         printf("%s", msg);
4
5     if (WIFEXITED(status)) {
6         printf("child exited, status=%d\n", WEXITSTATUS(status));
7     }
8     else if (WIFSIGNALED(status)) {
9         printf("child killed by signal %d (%s)",
10              WTERMSIG(status), strsignal(WTERMSIG(status)));
11         if (WCOREDUMP(status))
12             printf(" (core dumped)");
13         printf("\n");
14     }
15     else if (WIFSTOPPED(status)) {
16         printf("child stopped by signal %d (%s)\n",
17              WSTOPSIG(status), strsignal(WSTOPSIG(status)));
18     }
19     else if (WIFCONTINUED(status))
20         printf("child continued\n");
21 }
```

## An older wait API: *wait()*

---

```
#include <sys/wait.h>
pid_t wait(int *wstatus);
```

- The original “wait” API
- Equivalent to: `waitpid(-1, &wstatus, 0);`
- Still commonly used to handle the simple, common case:  
**wait for any child to terminate**

## Outline

---

6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: <code>fork()</code>	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The <code>SIGCHLD</code> signal	6-38
6.7	Executing programs: <code>execve()</code>	6-42

## Orphans

---

- An **orphan** is a process that lives longer than its parent
- Orphaned processes are **adopted by *init***
- ***init* waits for its adopted children** when they terminate
- After orphan is adopted, *getppid()* returns PID of *init*
  - Conventionally, *init* has PID 1
- On systems where the *init* system is *systemd*, then, depending on the configuration, things are different:
  - A helper process (PID  $\neq$  1) becomes parent of orphaned children
    - When run with the *--user* option, *systemd* organizes all processes in the user's session into a subtree with such a subreaper
  - See discussion of PR\_SET\_CHILD\_SUBREAPER in *prctl(2)*

[TLPI §26.2]

## Zombies

---

- Suppose a **child terminates before parent waits** for it
- Parent must still be able to collect status later
- $\Rightarrow$  Child becomes a **zombie**:
  - Most process resources are recycled
  - A process slot is retained
    - PID, status, and resource usage statistics
- Zombie is removed when parent does a “wait”

[TLPI §26.2]

## Creating a zombie: procexec/zombie.c

```
int main(int argc, char *argv[]) {
    int nzombies = (argc > 1) ? atoi(argv[1]) : 1;
    printf("Parent (PID %ld)\n", (long) getpid());

    for (int j = 0; j < nzombies; j++) {
        switch (fork()) {
            case -1:
                errExit("fork-%d", j);
            case 0:
                /* Child: exits to become zombie */
                printf("Child (PID %ld) exiting\n", (long) getpid());
                exit(EXIT_SUCCESS);
            default:
                /* Parent continues in loop */
                break;
        }
    }

    sleep(600); /* Children are zombies during this time */
    while (wait(NULL) > 0) /* Reap zombie children */
        continue;
    exit(EXIT_SUCCESS);
}
```

- Create one or more zombie child processes

## Creating a zombie: procexec/zombie.c

```
1 $ ./zombie &
2 [1] 23425
3 Parent (PID 23425)
4 Child (PID 23427) exiting
5 $ ps -C zombie
6   PID TTY          TIME CMD
7 23425 pts/1        00:00:00 zombie
8 23427 pts/1        00:00:00 zombie <defunct>
9 $ kill -KILL 23427
10 $ ps -C zombie
11   PID TTY          TIME CMD
12 23425 pts/1        00:00:00 zombie
13 23427 pts/1        00:00:00 zombie <defunct>
```

- **Zombies can't be killed** by signals!
  - (Since parent must still be able to “wait”)
  - Even silver bullets (SIGKILL) don't work

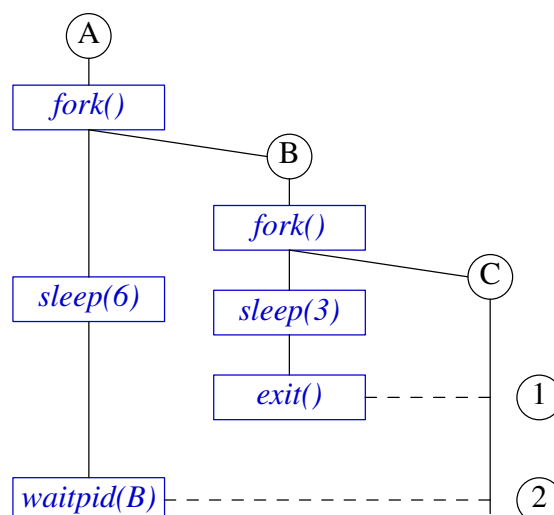


## Reap your zombies

- **Zombie may live for ever**, if parent fails to “wait” on it
  - Or until parent is killed, so zombie is adopted by *init*
- **Long-lived processes that create children must ensure that zombies are “reaped”** (“waited” for)
  - Shells, network servers, ...

## Exercise

- ① Suppose that we have three processes related as grandparent, parent, and child, and that the parent exits after a few seconds, but the grandparent does **not** immediately perform a *wait()* after the parent exits, with the result that the parent becomes a zombie, as in the following diagram.



## Exercise

---

When do you expect the child to be adopted by *init* (so that *getppid()* in the child returns 1): after the parent terminates or after the grandparent does a *wait()*? In other words, is the child adopted at point 1 or point 2 in the diagram? Write a program, **[(minimal) template: `procexec/ex.grandchild_zombie.c`]**, to verify the answer.

Note the following points:

- You will probably want to use calls to *sleep()* so that you can more easily observe the steps that occur during execution of the program. For example:
  - The grandchild could loop 10 times, displaying the value returned by *getppid()* and sleeping for 1 second on each loop iteration.
  - The parent could sleep for 3 seconds before terminating.
  - The grandparent could sleep for 6 seconds before calling *wait()*
- Depending on your distribution (e.g., if you have a *systemd*-based system where the `--user` flag is employed), you will see slightly different results from those described above. In particular, the orphaned child is reparented to a process other than PID 1.

## Outline

---

6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: <code>fork()</code>	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The <code>SIGCHLD</code> signal	6-38
6.7	Executing programs: <code>execve()</code>	6-42

# The SIGCHLD signal

---

- SIGCHLD is generated for a parent when a child terminates
- Ignored by default
- Catching SIGCHLD allows us to be asynchronously notified of child's termination
  - Can be more convenient than synchronous or nonblocking *waitpid()* calls
- Within SIGCHLD handler, we “wait” to reap zombie child

[TLPI §26.3]

## A SIGCHLD handler

---

```
1 void grimReaper(int sig) {  
2     int savedErrno = errno;  
3     while (waitpid(-1, NULL, WNOHANG) > 0)  
4         continue;  
5     errno = savedErrno;  
6 }
```

- Each *waitpid()* call reaps one terminated child
- while loop handles possibility that multiple children terminated while SIGCHLD was blocked
  - e.g., during earlier invocation of handler
- WNOHANG ensures handler does not block if there are no more terminated children
- Loop terminates when *waitpid()* returns:
  - 0, meaning no more *terminated* children
  - -1, probably with *errno == ECHILD*, meaning no more children
- Handler saves and restores *errno*, so that it is reentrant

## SIGCHLD for stopped and continued children

---

- SIGCHLD is also generated when a child stops or continues
- To prevent this, specify SA\_NOCLDSTOP in *sa\_flags* when establishing SIGCHLD handler with *sigaction()*

[TLPI §26.3.2]

## Outline

---

6	Process Lifecycle	6-1
6.1	Introduction	6-3
6.2	Creating a new process: <code>fork()</code>	6-6
6.3	Process termination	6-13
6.4	Monitoring child processes	6-19
6.5	Orphans and zombies	6-30
6.6	The SIGCHLD signal	6-38
6.7	Executing programs: <code>execve()</code>	6-42

## Executing a new program

---

`execve()` loads a new program into caller's memory

- Old program, stack, data, and heap are discarded
- After executing run-time start-up code, execution commences in new program's `main()`
- Various functions layered on top of `execve()`:
  - Provide variations on functionality of `execve()`
  - Collectively termed “`exec()`”
    - See `exec(3)` man page

[TLPI §27.1]

## Executing a new program with `execve()`

---

```
#include <unistd.h>
int execve(const char *pathname, char *const argv[],
           char *const envp[]);
```

- `execve()` loads program at `pathname` into caller's memory
- `pathname` is an absolute or relative pathname
- `argv` specifies command-line arguments for new program
  - Defines `argv` argument for `main()` in new program
  - NULL-terminated array of pointers to strings
- `argv[0]` is command name
  - Normally same as basename part of `pathname`
  - Program can vary its behavior, depending on value of `argv[0]`
    - `busybox`

## Executing a new program with `execve()`

```
#include <unistd.h>
int execve(const char *pathname, char *const argv[],
           char *const envp[]);
```

- *envp* specifies environment list for new program
  - Defines *environ* in new program
  - NULL-terminated array of pointers to strings
- Successful `execve()` does not return
- If `execve()` returns, it failed; no need to check return value:

```
execve(pathname, argv, envp);
printf("execve() failed\n");
```

## Example: `procexec/exec_status.c`

```
./exec_status command [args...]
```

- Create a child process
- Child executes *command* with supplied command-line arguments
- Parent waits for child to exit, and reports wait status

## Example: procexec/exec\_status.c

```
1 extern char **environ;
2 int main(int argc, char *argv[]) {
3     pid_t childPid, wpid;
4     int wstatus;
5     ...
6     switch (childPid = fork()) {
7         case -1: errExit("fork");
8
9         case 0: /* Child */
10             printf("PID of child: %ld\n", (long) getpid());
11             execve(argv[1], &argv[1], environ);
12             errExit("execve");
13
14         default: /* Parent */
15             wpid = waitpid(childPid, &wstatus, 0);
16             if (wpid == -1) errExit("waitpid");
17             printf("Wait returned PID %ld\n", (long) wpid);
18             printWaitStatus(" ", wstatus);
19     }
20     exit(EXIT_SUCCESS);
21 }
```

## Example: procexec/exec\_status.c

```
1 $ ./exec_status /bin/date
2 PID of child: 4703
3 Thu Oct 24 13:48:44 NZDT 2013
4 Wait returned PID 4703
5     child exited, status=0
6 $ ./exec_status /bin/sleep 60 &
7 [1] 4771
8 PID of child: 4773
9 $ kill 4773
10 Wait returned PID 4773
11     child killed by signal 15 (Terminated)
12 [1]+  Done          ./exec_status /bin/sleep 60
```

## Exercise

---

- ① Write a simple shell program. The program should loop, continuously reading shell commands from standard input. Each input line consists of a set of white-space delimited words that are a command and its arguments. Each command should be executed in a new child process (*fork()*) using *execve()*. The parent process (the “shell”) should wait on each child and display its wait status (you can use the supplied *printWaitStatus()* function).

**[template: procexec/ex.simple\_shell.c]**

Some hints:

- The space-delimited words in the input line need to be broken down into a set of null-terminated strings pointed to by an *argv*-style array, and that array must end with a NULL pointer. The *strtok(3)* library function simplifies this task. (This task is already performed by code in the template.)
- Because *execve()* is used, you will need to type the full pathname when entering commands to your shell

## Exercise

---

- ② Write a program, *procexec/exec\_self\_pid.c*, that verifies that an *exec* does not change a process's PID
- The program should perform the following steps:
    - Print the process's PID.
    - If *argc* is 2, the program exits.
    - Otherwise, the program uses *exec()* to re-execute itself with an additional command-line argument (any string), so that *argc* will be 2.
  - Test the program by running it with no command-line arguments (i.e., *argc* is 1).



## Exercise

- ③ Write a program ([**template:** `procexec/ex.make_link.c`]) that takes two arguments:

```
make_link target linkpath
```

If invoked with the name *slink*, it creates a symbolic link (*symlink()*) using these pathnames, otherwise it creates a hard link (*link()*). After compiling, create two hard links to the executable, with the names *hlink* and *slink*. Verify that when run with the name *hlink*, the program creates hard links, while when run with the name *slink*, it creates symbolic links.

Hint:

- You will find the *basename()* and *strcmp()* functions useful when inspecting the program name in *argv[0]*.

## The `exec()` library functions

```
#include <unistd.h>
int execl(const char *pathname, const char *arg, ...
          /* , (char *) NULL, char *const envp[] */);
int execlp(const char *filename, const char *arg, ...
          /* , (char *) NULL */);
int execvp(const char *filename, char *const argv[]);
int execv(const char *pathname, char *const argv[]);
int execl(const char *pathname, const char *arg, ...
          /* , (char *) NULL */);
int execvpe(const char *filename, const char *argv[],
            char *const envp[]);
```

- Variations on theme of `execve()`
- Like `execve()`, the `exec()` functions return only if they fail
- `execvpe()` is Linux-specific (define `_GNU_SOURCE`)

# The `exec()` library functions

Vary theme of `execve()` with 2 choices in each of 3 dimensions:

- How are command-line arguments of new program specified?
- How is the executable specified?
- How is environment of new program specified?

Final letters in name of each function are clue about behavior

Function	Specification of arguments (v, l)	Specification of executable file (-, p)	Source of environment (e, -)
<code>execve()</code>	array	pathname	<i>envp</i> argument
<code>execl()</code>	list	pathname	<i>envp</i> argument
<code>execlp()</code>	list	filename + PATH	caller's <i>environ</i>
<code>execvp()</code>	array	filename + PATH	caller's <i>environ</i>
<code>execv()</code>	array	pathname	caller's <i>environ</i>
<code>execl()</code>	list	pathname	caller's <i>environ</i>
<code>execvpe()</code>	array	filename + PATH	<i>envp</i> argument

## Notes

Notes

Notes

Linux System Programming Essentials

# System Call Tracing with strace

Michael Kerrisk, man7.org © 2021

mtk@man7.org

November 2021

## Outline

---

7	System Call Tracing with strace	7-1
7.1	Getting started	7-3
7.2	Tracing child processes	7-10
7.3	Filtering strace output	7-14
7.4	System call tampering	7-20
7.5	Further strace options	7-26

# Outline

---

7	System Call Tracing with strace	7-1
7.1	Getting started	7-3
7.2	Tracing child processes	7-10
7.3	Filtering strace output	7-14
7.4	System call tampering	7-20
7.5	Further strace options	7-26

## *strace(1)*

---

- A tool to trace system calls made by a user-space process
  - Implemented via *ptrace(2)*
- Or: a debugging tool for tracing **complete conversation between application and kernel**
  - Application source code is not required
- Answer questions like:
  - What system calls are employed by application?
  - Which files does application touch?
  - What arguments are being passed to each system call?
  - Which system calls are failing, and why (*errno*)?

## *strace(1)*

- Log information is provided in **symbolic form**
  - **System call names** are shown
  - We see **signal names** (not numbers)
  - **Strings** printed as characters (up to 32 bytes, by default)
  - **Bit-mask arguments displayed symbolically**, using corresponding bit flag names ORed together
  - **Structures** displayed with **labeled fields**
  - **errno values** displayed symbolically + matching error text
  - “large” arguments and structures are abbreviated by default

```
fstat(3, {st_dev=makedev(0x8, 0x5), st_ino=407279,
st_mode=S_IFREG|0755, st_nlink=1, st_uid=0, st_gid=0,
st_blksize=4096, st_blocks=80, st_size=36960, st_atime=1625615479
/* 2021-07-07T01:51:19.795021222+0200 */ , st_atime_nsec=795021222,
st_mtime=1613345143 /* 2021-02-15T00:25:43+0100 */ , st_mtime_nsec=0,
st_ctime=1616161103 /* 2021-03-19T14:38:23.816838407+0100 */ ,
st_ctime_nsec=816838407}) = 0
open("/lib64/liblzma.so.5", 0_RDONLY|O_CLOEXEC) = 3
```

## Simple usage: tracing a command at the command line

- A very simple C program:

```
int main(int argc, char *argv[]) {
#define STR "Hello world\n"
    write(STDOUT_FILENO, STR, strlen(STR));
    exit(EXIT_SUCCESS);
}
```

- Run *strace(1)*, directing logging output (*-o*) to a file:

```
$ strace -o strace.log ./hello_world
Hello world
```

- (By default, trace output goes to standard error)
- ⚠ On some systems, may first need to ensure *ptrace\_scope* file has value 0 or 1:

```
$ sudo sh -c 'echo 0 > /proc/sys/kernel/yama/ptrace_scope'
```

- Yama LSM disables *ptrace(2)* to prevent attack escalation; see *ptrace(2)* man page

## Simple usage: tracing a command at the command line

```
$ cat strace.log
execve("./hello_world", ["/hello_world"], [/* 110 vars */]) = 0
...
access("/etc/ld.so.preload", R_OK) = -1 ENOENT
(No such file or directory)
open("/etc/ld.so.cache", O_RDONLY|O_CLOEXEC) = 3
fstat(3, {st_mode=S_IFREG|0644, st_size=160311, ...}) = 0
mmap(NULL, 160311, PROT_READ, MAP_PRIVATE, 3, 0) = 0x7fa5ecfc0000
close(3) = 0
open("/lib64/libc.so.6", O_RDONLY|O_CLOEXEC) = 3
...
write(1, "Hello world\n", 12) = 12
exit_group(0) = ?
+++ exited with 0 +++
```

- Even simple programs make lots of system calls!
  - 25 in this case (many have been edited from above output)
- Most output in this trace relates to finding and loading shared libraries
  - First call (`execve()`) was used by shell to load our program
  - Only last two system calls were made by our program

## Simple usage: tracing a command at the command line

```
$ cat strace.log
execve("./hello_world", ["/hello_world"], [/* 110 vars */]) = 0
...
access("/etc/ld.so.preload", R_OK) = -1 ENOENT
(No such file or directory)
open("/etc/ld.so.cache", O_RDONLY|O_CLOEXEC) = 3
fstat(3, {st_mode=S_IFREG|0644, st_size=160311, ...}) = 0
mmap(NULL, 160311, PROT_READ, MAP_PRIVATE, 3, 0) = 0x7fa5ecfc0000
close(3) = 0
open("/lib64/libc.so.6", O_RDONLY|O_CLOEXEC) = 3
...
write(1, "Hello world\n", 12) = 12
exit_group(0) = ?
+++ exited with 0 +++
```

For each system call, we see:

- Name of system call
- Values passed in/returned via arguments
- System call return value
- Symbolic *errno* value (+ explanatory text) on syscall failures

## A gotcha...

---

- The last call in our program was:

```
exit(EXIT_SUCCESS);
```

- But *strace* showed us:

```
exit_group(0)                                = ?
```

- Some detective work:

- We “know” *exit(3)* is a library function that calls *\_exit(2)*
- But where did *exit\_group()* come from?
- *\_exit(2)* man page tells us:

```
$ man 2 _exit
...
C library/kernel differences
In glibc up to version 2.3, the _exit() wrapper function
invoked the kernel system call of the same name. Since
glibc 2.3, the wrapper function invokes exit_group(2),
in order to terminate all of the threads in a process.
```

- $\Rightarrow$  may need to dig deeper to understand *strace(1)* output

## Outline

---

7	System Call Tracing with <i>strace</i>	7-1
7.1	Getting started	7-3
7.2	Tracing child processes	7-10
7.3	Filtering <i>strace</i> output	7-14
7.4	System call tampering	7-20
7.5	Further <i>strace</i> options	7-26



## Tracing child processes

- By default, *strace* does not trace children of traced process
- *-f* option causes children to be traced
  - Each trace line is prefixed by PID
  - In a program that employs POSIX threads, each line shows kernel thread ID (*gettid()*)

## Tracing child processes: *strace/fork\_exec.c*

```
1 int main(int argc, char *argv[]) {
2     pid_t childPid;
3     char *newEnv[] = {"ONE=1", "TWO=2", NULL};
4
5     printf("PID of parent: %ld\n", (long) getpid());
6     childPid = fork();
7     if (childPid == 0) {          /* Child */
8         printf("PID of child:  %ld\n", (long) getpid());
9         if (argc > 1) {
10             execve(argv[1], &argv[1], newEnv);
11             errExit("execve");
12         }
13         exit(EXIT_SUCCESS);
14     }
15     wait(NULL);                  /* Parent waits for child */
16     exit(EXIT_SUCCESS);
17 }
```

```
$ strace -f -o strace.log ./fork_exec
PID of parent: 1939
PID of child:  1940
```

# Tracing child processes: `strace/fork_exec.c`

```
$ cat strace.log
1939 execve("./fork_exec", ["/fork_exec"], [/* 110 vars */]) = 0
...
1939 clone(child_stack=0, flags=CLONE_CHILD_CLEARTID|
    CLONE_CHILD_SETTID|SIGCHLD, child_tidptr=0x7fe484b2ea10) = 1940
1939 wait4(-1, <unfinished ...>
1940 write(1, "PID of child:  1940\n", 21) = 21
1940 exit_group(0) = ?
1940 +++ exited with 0 +++
1939 <... wait4 resumed> NULL, 0, NULL) = 1940
1939 --- SIGCHLD {si_signo=SIGCHLD, si_code=CLD_EXITED,
    si_pid=1940, si_uid=1000, si_status=0, si_utime=0,
    si_stime=0} ---
1939 exit_group(0) = ?
1939 +++ exited with 0 +++
```

- Each line of trace output is prefixed with corresponding PID
- Inside glibc, `fork()` is actually a wrapper that calls `clone(2)`
- `wait()` is a wrapper that calls `wait4(2)`
- We see two lines of output for `wait4()` because call blocks and then resumes
- `strace` shows us that parent received a SIGCHLD signal

## Outline

7	System Call Tracing with <code>strace</code>	7-1
7.1	Getting started	7-3
7.2	Tracing child processes	7-10
7.3	Filtering <code>strace</code> output	7-14
7.4	System call tampering	7-20
7.5	Further <code>strace</code> options	7-26

## Selecting system calls to be traced

- `strace -e` can be used to select system calls to be traced
- `-e trace=<syscall>[, <syscall>...]`
  - Specify system call(s) that should be traced
  - Other system calls are ignored

```
$ strace -o strace.log -e trace=open,close ls
```

- `-e trace=!<syscall>[, <syscall>...]`
  - **Exclude** specified system call(s) from tracing
    - Some applications do bizarre things (e.g., calling `gettimeofday()` 1000s of times/sec.)
  - `!` needs to be quoted to avoid shell interpretation
- `-e trace=/<regex>`
  - Trace syscalls whose names match regular expression
    - April 2017; expression will probably need to be quoted...

## Selecting system calls by category

- `-e trace=<syscall-category>` trace a category of syscalls
- Categories include:
  - `%file`: trace all syscalls that take a filename as argument
    - `open()`, `stat()`, `truncate()`, `chmod()`, `setxattr()`, `link()`...
  - `%desc`: trace file-descriptor-related syscalls
    - `read()`, `write()`, `open()`, `close()`, `fsetxattr()`, `poll()`, `select()`, `pipe()`, `fcntl()`, `epoll_create()`, `epoll_wait()`...
  - `%process`: trace process management syscalls
    - `fork()`, `clone()`, `exit_group()`, `execve()`, `wait4()`, `unshare()`...
  - `%network`: trace network-related syscalls
    - `socket()`, `bind()`, `listen()`, `connect()`, `sendmsg()`...
  - `%signal`: trace signal-related syscalls
    - `kill()`, `rt_sigaction()`, `rt_sigprocmask()`, `rt_sigqueueinfo()`...
  - `%memory`: trace memory-mapping-related syscalls
    - `mmap()`, `mprotect()`, `mlock()`...

## Filtering signals

- *strace -e signal=set*
  - Trace only specified set of signals
  - “sig” prefix in names is optional; following are equivalent:

```
$ strace -e signal=sigio,sigint ls > /dev/null
$ strace -e signal=io,int ls > /dev/null
```

- *strace -e signal=!set*
  - Exclude specified signals from tracing

## Filtering by pathname

- *strace -P pathname*: trace only system calls that access file at *pathname*
  - Specify multiple *-P* options to trace multiple paths
- Example:

```
$ strace -o strace.log -P /lib64/libc.so.6 ls > /dev/null
Requested path '/lib64/libc.so.6' resolved into
'/usr/lib64/libc-2.18.so'
$ cat strace.log
open("/lib64/libc.so.6", 0_RDONLY|0_CLOEXEC) = 3
read(3, "\177ELF\2\1\1\3\0\0\0\0\0\0\0\0\3\0>\0\1\0\0\0p\36
\2\0\0\0\0\0"... , 832) = 832
fstat(3, {st_mode=S_IFREG|0755, st_size=2093096, ...}) = 0
mmap(NULL, 3920480, PROT_READ|PROT_EXEC,
      MAP_PRIVATE|MAP_DENYWRITE, 3, 0) = 0x7f8511fa3000
mmap(0x7f8512356000, 24576, PROT_READ|PROT_WRITE,
      MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3, 0x1b3000)
      = 0x7f8512356000
close(3)                                = 0
+++ exited with 0 +++
```

- *strace* noticed that the specified file was opened on FD 3, and also traced operations on that FD

# Mapping file descriptors to pathnames

---

- `-y` option causes *strace* to display pathnames corresponding to each file descriptor
  - Useful info is also displayed for other types of file descriptors, such as pipes and sockets

```
$ strace -y cat greet
...
openat(AT_FDCWD, "greet", O_RDONLY) = 3</home/mtk/greet>
fstat(3</home/mtk/greet>, {st_mode=S_IFREG|0644, ...
read(3</home/mtk/greet>, "hello world\n", 131072) = 12
write(1</dev/pts/11>, "hello world\n", 12) = 12
read(3</home/mtk/greet>, "", 131072) = 0
close(3</home/mtk/greet>) = 0
...
```

- `-yy` is as for `-y` but shows additional protocol-specific info for sockets

```
write(3<TCP: [10.0.20.135:33522->213.131.240.174:80]>,
"GET / HTTP/1.1\r\nUser-Agent: Wget"... , 135) = 135
read(3<TCP: [10.0.20.135:33522->213.131.240.174:80]>,
"HTTP/1.1 200 OK\r\nDate: Thu, 19 J"... , 253) = 253
```

## Outline

---

7	System Call Tracing with <i>strace</i>	7-1
7.1	Getting started	7-3
7.2	Tracing child processes	7-10
7.3	Filtering <i>strace</i> output	7-14
7.4	System call tampering	7-20
7.5	Further <i>strace</i> options	7-26

## System call tampering

---

- *strace* can be used to **modify** behavior of selected syscall(s)
  - Initial feature implementation completed in early 2017
- Various possible effects:
  - Inject delay before/after syscall
  - Generate a signal on syscall
  - Bypass execution of syscall, making it return a “success” value or fail with specified value in *errno* (error injection)
  - (Limited) ability to choose which invocation of syscall will be modified

## *strace* -e inject options

---

- Syntax: *strace* -e inject=<syscall-set>[:<option>]...
  - *syscall-set* is set of syscalls whose behavior will be modified
- :error=*errnum*: syscall is not executed; returns failure status with *errno* set as specified
- :retval=*value*: syscall is not executed; returns specified “success” value
  - Can't specify both :retval and :error together

## strace -e inject options

- `:signal=sig`: deliver specified signal on entry to syscall
- `:delay_enter=usecs, :delay_exit=usecs`: delay for *usecs* microseconds on entry to/return from syscall
- `:when=expr`: specify which invocation(s) to tamper with
  - `:when=N`: tamper with invocation *N*
  - `:when=N+`: tamper starting at *N*th invocation
  - `:when=N+S`: tamper with invocation *N*, and then every *S* invocations
  - Range of *N* and *S* is 1..65535

## Example

```
$ strace -y -e close \
    -e inject=close:error=22:when=3 /bin/ls > d
close(3</etc/ld.so.cache>) = 0
close(3</usr/lib64/libselinux.so.1>) = 0
close(3</usr/lib64/libcap.so.2.25>) = -1 EINVAL
(Invalid argument) (INJECTED)
close(3</usr/lib64/libcap.so.2.25>) = 0
/bin/ls: error while loading shared libraries: libcap.so.2:
cannot close file descriptor: Invalid argument
+++ exited with 127 +++
```

- Use `-y` to show pathnames corresponding to file descriptors
- Inject error 22 (EINVAL) on third call to `close()`
- Third `close()` was not executed; an error return was injected
  - (After that, *ls* got sad)

## Using system call tampering for error injection

---

- Success-injection example: make `unlinkat()` succeed, without deleting temporary file that would have been deleted
- Error-injection use case: quick and simple black-box testing
  - Does application fail gracefully when encountering unexpected error?
- But there are alternatives for black-box testing:
  - Preloaded library with interposing wrapper function that spoofs a failure (without calling “real” function)
    - Can be more flexible
    - But can’t be used with set-UID/set-GID programs
  - Seccomp (secure computing)
    - Generalized facility to block execution of system calls based on system call number and argument values
    - More powerful, but can’t, for example cause Nth call to fail

## Outline

---

7	System Call Tracing with strace	7-1
7.1	Getting started	7-3
7.2	Tracing child processes	7-10
7.3	Filtering strace output	7-14
7.4	System call tampering	7-20
7.5	Further strace options	7-26



## Obtaining a system call summary

- *strace -c* counts time, calls, and errors for each system call and reports a summary on program exit

```
$ strace -c who > /dev/null
```

% time	seconds	usecs/call	calls	errors	syscall
21.77	0.000648	9	72		alarm
14.42	0.000429	9	48		rt_sigaction
13.34	0.000397	8	48		fcntl
8.84	0.000263	5	48		read
7.29	0.000217	13	17	2	kill
6.79	0.000202	6	33	1	stat
5.41	0.000161	5	31		mmap
4.44	0.000132	4	31	6	open
2.89	0.000086	3	29		close
2.86	0.000085	43	2		socket
2.82	0.000084	42	2	2	connect
...					
100.00	0.002976		442	13	total

- Treat time measurements as indicative only, since *strace* adds overhead to each syscall

## Tracing live processes

- *-p PID*: **trace running process** with specified PID
  - Type *Control-C* to cease tracing
  - To **trace multiple processes**, specify *-p* multiple times
  - Can trace only processes you own
  - ⚠️ ⚠️ tracing a process can **heavily affect performance**
    - E.g., up to two orders of magnitude slow-down in syscalls
    - ⚠️ Think twice before using in a production environment
- *-p PID -f*: will **trace all threads** in specified process

## Further *strace* options

---

- `-k`: print a stack trace after each traced syscall
- `sudo strace -u <username> prog`: run program with UID and GIDs of specified user
  - Useful when tracing privileged programs, such as set-UID-*root* programs
    - Normally, privileged programs are **not** run with privilege when executed under *strace*

## Further *strace* options

---

- `-v`: don't abbreviate arguments (structures, etc.)
  - Output can be quite verbose...
- `-s strsize`: maximum number of bytes to display for strings
  - Default is 32 characters
  - Pathnames are always printed in full
- Various options show start time or duration of system calls
  - `-t`, `-tt`: prefix each trace line with wall-clock time
    - `-tt` also adds microseconds
  - `-T`: show time spent in syscall
    - But treat as indications only, since *strace* causes overhead on syscalls

Notes

Notes