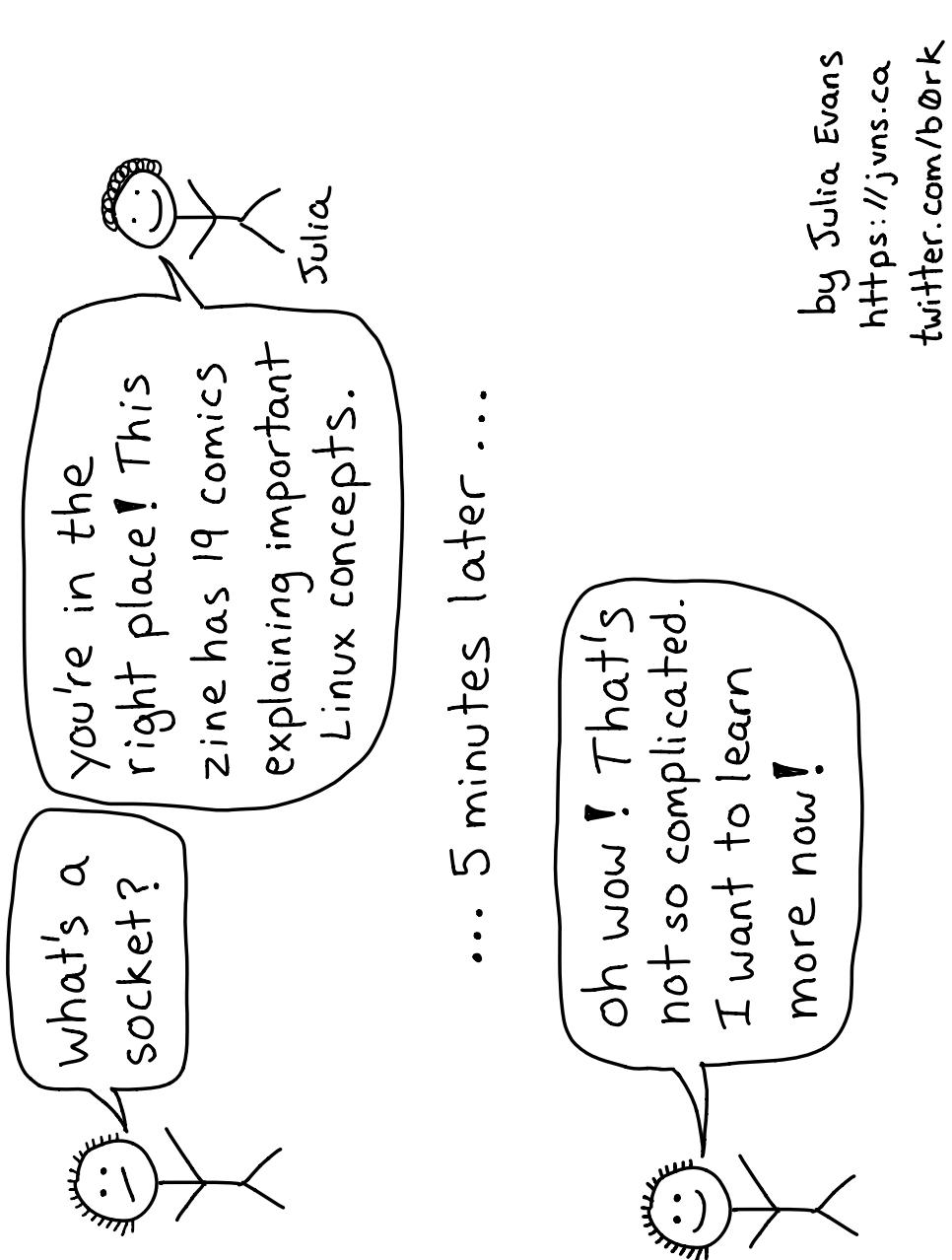


By Julia Evans

love this?
find more awesome zines at
→ wizardzines.com ←



by Julia Evans
<https://jvns.ca>
twitter.com/bork

THE LINUX PROGRAMMING INTERFACE

MICHAEL KERRISK

Want to learn more?
I highly recommend
this book:

Every chapter is a readable,
short (usually 10-20 pages)
explanation of a Linux system.

I used it as a reference
constantly when writing
this zine.

I ♡ it because even though
it's huge and comprehensive
(1500 pages!), the chapters
are short and self-contained
and it's very easy to pick it
up and learn something.



man page sections 22

man pages are split up
into 8 sections

① ② ③ ④ ⑤ ⑥ ⑦ ⑧

\$ man 2 read

means "get me the man page
for `read` from section 2".

There's both

→ a program called "read"
→ and a system call called "read"

so

\$ man 1 read

gives you a different man page from

\$ man 2 read

If you don't specify a section, man will
look through all the sections & show
the first one it finds.

man page sections

① programs

② system calls

\$ man grep

\$ man sendfile

\$ man ls

\$ man ptrace

③ C functions

④ devices

\$ man printf

\$ man null

\$ man fopen

for /dev/null docs

⑤ file formats

⑥ games

\$ man sudoers

\$ man s!

\$ man proc

files in /proc!

⑦ miscellaneous

⑧ sysadmin programs

explains concepts!

\$ man 7 pipe

\$ man apt

\$ man 7 symlink

\$ man chroot

♥ Table of contents ♥

unix permissions...4	sockets.....10	virtual memory...17
/proc.....5	unix domain sockets.....11	shared libraries...18
system calls.....6	processes.....12	copy on write....19
signals.....7	threads.....13	page faults.....20
file descriptors.8	floating point....14	mmap.....21
pipes.....9	file buffering.....15	man page sections.....22
memory allocation.....16		

Unix permissions

4

There are 3 things you can do to a file

`ls -l file.txt` shows you permissions.
Here's how to interpret the output:

→ **read** → **write** → **execute**

`rw-` → `rwx` → `r--` → bork staff

bork (user) **staff** (group)
can can
read & write read & write

File permissions are 12 bits

setuid ↗ user group all
000 110 110 100
sticky ↗ **rwx** ↗ **rwx** ↗ **rwx**
For files:
r = can read
w = can write
x = can execute
For directories, it's approximately:
r = can list files
w = can create files
x = can cd into & access files

File permissions are 12 bits
110 in binary is 6
So `rw-` `r--` `r--`
= 110 100 100
= 6 4 4
chmod 644 file.txt
means change the permissions to:
`rwx - r-- r--`
Simple!

setuid affects executables
`$ ls -l /bin/ping`
`rws -r-x r-x root root`
this means ping always runs as root
setgid does 3 different unrelated things for executables, directories, and regular files.
unix! it's a long story! why?

mmap

21

what's mmap for?

I want to work with a VERY LARGE FILE but it won't fit in memory

You could try mmap!
(mmap = "memory map")

load files lazily with mmap
When you mmap a file, it gets mapped into your program's memory.
2TB file ↗ virtual memory but nothing is ACTUALLY read into RAM until you try to access the memory.
(how it works: page faults!)

sharing big files with mmap

we all want to read the same file!
no problem! mmap

Even if 10 processes mmap a file, it will only be read into memory once.

how to mmap in Python
`import mmap`
`f = open("HUGE.txt")`
`mm = mmap.mmap(f.fileno(), 0)`
this won't read the file from disk!
Finishes ~instantly.
`print(mm[-1000:1])`
this will read only the last 1000 bytes!

anonymous memory maps

→ not from a file (memory set to 0 by default)
→ with MAP_SHARED, you can use them to share memory with a subprocess!

dynamic linking
uses mmap
I need to use libc.so.6 C standard library
you too eh? no problem!
I always mmap, so that file is probably dynamic loaded into memory already.

page faults

every Linux process has a page table

* page table *

virtual memory address	physical memory address
0x19723000	0x1422000
0x19724000	0x1423000
0x1524000	not in memory
0x1844000	0x4a000 read only

"not resident in memory" usually means the data is on disk!

virtual memory



Having some virtual memory that is actually on disk is how swap and mmap work.

some pages are marked as either

* read only

not resident in memory when you try to access a page that's marked "not resident in memory," it triggers a !page fault!

what happens during a page fault?

- your program stops running
- Linux kernel code to handle the page fault runs
- Linux running

I'll fix the problem and let your program keep running

an amazing directory: /proc 5

/proc/PID/cmdline

command line arguments the process was started with

/proc/PID/exe

symlink to the process's binary **magic**: works even if the binary has been deleted!

/proc/PID/environ

all of the process's environment variables

/proc/PID/status

Is the program running or asleep? How much memory is it using? And much more!

/proc/PID/fd

Directory with every file the process has open!

Run `ls -l /proc/42/fd` to see the list of files for process 42.

These symlinks are also magic & you can use them to recover deleted files ❤️

/proc/PID/stack

The kernel's current stack for the process. Useful if it's stuck in a system call.

/proc/PID/maps

List of process's memory maps. Shared libraries, heap, anonymous maps, etc.

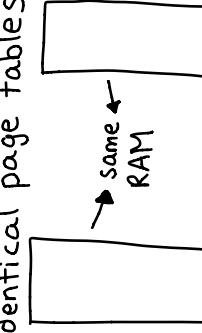
and ;;;;
Look at
man proc

for more information!

System calls

<p>The Linux kernel has code to do a lot of things</p> <ul style="list-style-type: none">read from a hard drivemake network connectionscreate new processkill processchange file permissionskeyboard drivers	<p>your program <u>doesn't</u> know how to do those things</p> <ul style="list-style-type: none">TCP? dude I have no idea how that works.NO, I do not know how the ext4 filesystem is implemented. I just want to read some files!	<p>every program uses system calls</p> <ul style="list-style-type: none">I use the 'open' syscall to open filesme too!Java programme three!C program
<p>programs ask Linux to do work for them using <u>system calls</u>:</p> <ul style="list-style-type: none">please write to this fileswitch to running kernel code	<p>done! I wrote 1097 bytes!</p> <ul style="list-style-type: none">Linux<program resumes>	<p>and every system call has a number (e.g. chmod is #90 on x86-64)</p> <p>So what's actually going on when you change a file's permissions is:</p> <ul style="list-style-type: none">run syscall #90 program with these argumentsok!Linux
		<p>you can see which system calls a program is using with <u>strace</u>:</p> <ul style="list-style-type: none">\$ strace ls /tmpwill show you every system call 'ls' uses!it's really fun!strace has high overhead so don't run it on your production database

<p>programs ask Linux to do work for them using <u>system calls</u>:</p> <ul style="list-style-type: none">please write to this fileswitch to running kernel code	<p>done! I wrote 1097 bytes!</p> <ul style="list-style-type: none">Linux<program resumes>	<p>you can see which system calls a program is using with <u>strace</u>:</p> <ul style="list-style-type: none">\$ strace ls /tmpwill show you every system call 'ls' uses!it's really fun!strace has high overhead so don't run it on your production database
--	--	---

<h2>Copy on write</h2> <p>On Linux, you start new processes using the <u>fork()</u> or <u>clone()</u> system call.</p> <p>calling fork creates a child process that's a copy of the caller</p> <ul style="list-style-type: none">parentchild	<p>copying all that memory every time we fork would be <u>slow</u> and a <u>waste of RAM</u></p> <p>often processes call <u>exec</u> right after <u>fork</u>, which means they don't use the parent process's memory basically at all!</p>	<p>so Linux lets them share physical RAM and only copies the memory when one of them tries to <u>write</u></p> <ul style="list-style-type: none">I'd like to change that memoryokay! I'll make you your own copy!Linux
	<p>Linux does this by giving both the processes identical page tables.</p> 	<p>but it marks every page as <u>read only</u>.</p>

<p>6</p>

Shared libraries

18

Most programs on Linux use a bunch of C libraries. Some popular libraries:

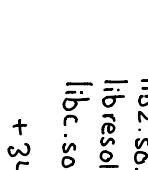
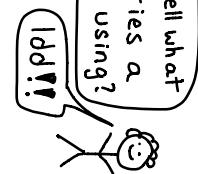
openssl
(for SSL!)

sqlite
(embedded db!)

libpcre
(regular
expressions!)

zlib
(gzip!)

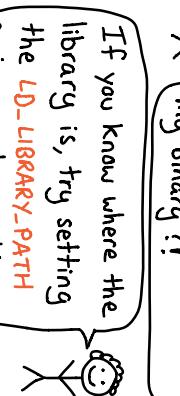
libstdc++
(C++ standard library!)

 how can I tell what shared libraries a program is using?
 lib!!!

```
$ ldd /usr/bin/curl
libz.so.1 => /lib/x86_64...
libresolv.so.2 => ...
libc.so.6 => ...
+ 34 more !!
```

There are 2 ways to use any library:
 ① Link it into your binary
 big binary with lots of things!
 [your code] **z lib** **sqlite**

② Use separate shared libraries
 [your code] **z lib** **sqlite** ← all different files

 I got a "library not found" error when running my binary !?
 If you know where the library is, try setting the **LD_LIBRARY_PATH** environment variable
 oo LD-LIBRARY-PATH tells me where to look!

① DT_RPATH in your executable
 ② LD-LIBRARY-PATH
 ③ DT_RUNPATH in executable
 ④ /etc/ld.so.cache (run ldconfig -p to see contents)
 ⑤ /lib, /usr/lib

Programs like this
 [your code] **z lib** **sqlite**
 are called "statically linked"
 and programs like this
 [your code] **z lib** **sqlite**
 are called "dynamically linked"

Signals

7

If you've ever used kill
you've used signals



Every signal has a default action, which is one of:

ignore

kill process

kill process AND make core dump file

stop process

resume process

the Linux kernel sends processes signals in lots of situations

your child terminated

that pipe is closed

illegal instruction

the timer you set expired

Segmentation fault

Your program can set custom handlers for almost any signal

SIGTERM (terminate)

clean up and then exit!

process

SIGSTOP & SIGKILL can't be ignored

got → **xx**

sigkilled **xx**

you can send signals yourself with the kill system call or command

SIGINT Ctrl-C

SIGTERM Kill

SIGKILL Kill -9

} various levels of "die"

SIGHUP kill -HUP

often interpreted as "reload config", e.g. by nginx

signals can be hard to handle correctly since they can happen at ANY time

oo handling a signal process

SURPRISE! another signal!

got → **xx**

file descriptors

8

Unix systems use integers to track open files

`ls -l /proc/PID/fd` shows you a process's open files

`$ ls -l /proc/4242/fd` → PID we're interested in

FD NAME	0	/dev/pts/pty1
	1	/dev/pts/pty1
	2	pipe:29174
	3	/home/bark/awesome.txt
	5	/tmp/

↑ FD is for file descriptor

process → kernel

Okay! That's file #7 for you.

these integers are called **file descriptors**

When you read or write to a file/pipe/network connection you do that using a file descriptor

connect to google.com → fd 5 → OS

write GET / HTTP/1.1 to fd #5 → done!

When you read from a file descriptor under the hood:

Python:

```
f = open("file.txt")
f.read_lines()
```

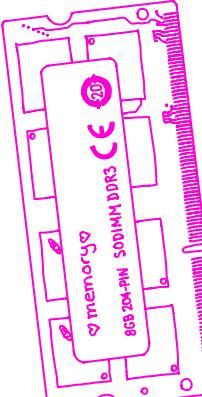
Behind the scenes:

open file.txt → Python program → read from file #4 → OS → ok! fd is 4 → Here are the contents!

virtual memory

17

your computer has physical memory



Linux keeps a mapping from virtual memory pages to physical memory pages called the **page table**

a "page" is a 4KB sometimes bigger chunk of memory

every time you switch which process is running, Linux needs to switch the page table

CPU → I'm accessing 0x21000

MMU → I'll look that up in the **page table** and then access the right physical address unit, hardware

here's the address of process 2950's page table

Linux → thanks, I'll use that now!

PID	virtual addr	physical addr
1971	0x20000	0x192000
2310	0x20000	0x228000
2310	0x21000	0x9788000

memory allocation

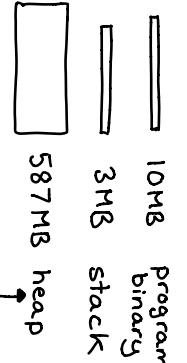
16

Your program has memory

(malloc)

is responsible for 2 things.

THING 1: keep track of what memory is used / free.



the heap is what your allocator manages

your memory allocator's interface

`malloc (size_t size)`

allocate `size` bytes of

memory & return a pointer to it.

`free (void* pointer)`

mark the memory as unused (and maybe give back to the OS).

`realloc(void* pointer, size_t size)`

ask for more / less memory for `pointer`.

`calloc(size_t members, size_t size)`

allocate array + initialize to 0.

your memory allocator (`malloc`) is responsible for more memory!

THING 2: Ask the OS

for more memory!



`malloc`

(here you go!) OS

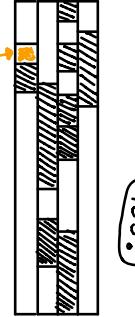
malloc tries to fill in unused space when you ask for memory

`malloc`

can I have 512 bytes of memory?

YES!

malloc



can I have 60MB more?

malloc

(here you go!) OS

malloc isn't `:magic:`!

it's just a function!

`malloc`

you can always:

→ use a different `malloc` library like `jemalloc`

or `tcmalloc` (easy!)

→ implement your own `malloc` (harder)

pipes

9

Sometimes you want to send the output of one process to the input of another

`$ ls | wc -l`

53
53 files!

Linux creates a buffer for each pipe.

`ls` → buffer → `wc`

`IN` [data waiting to be read] `out`

If data gets written to the pipe faster than it's read, the buffer will fill up. `IN` [██████████] `out`

When the buffer is full, writes to `IN` will block (wait) until the reader reads. This is normal & ok!)

a pipe is a pair of 2 magical file descriptors

① pipe input `IN` → `OUT`

② pipe output

`ls` → `IN` → `OUT` → `wc` → `stdout`

`stdin`

what if your target process dies?

`ls` → `IN` → `OUT` → `wc`

`IN` [xx] `OUT`

If `wc` dies, the pipe will close and `ls` will be sent SIGPIPE. By default, SIGPIPE terminates your process.

named pipes

`$ mkfifo my-pipe`

This lets 2 unrelated processes communicate through a pipe!

`f=open("./my-pipe")`

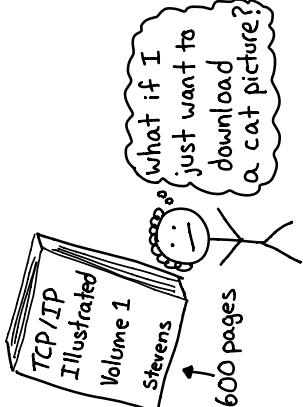
`f.write("hi\n")`

`f=open("./my-pipe")`
`f.readline() ← "hi!"`

Sockets

10

networking protocols are complicated



Unix systems have an API called the "socket API" that makes it easier to make network connections

You don't need to know how TCP works.
Unix I'll take care of it!

here's what getting a cat picture with the socket API looks like:

- ① Create a socket
`fd = socket(AF_INET, SOCK_STREAM ...)`
- ② Connect to an IP / port
`connect(fd, (fd, 12.13.14.15:80))`
- ③ Make a request
`write(fd, "GET /cat.png HTTP/1.1 ...")`
- ④ Read the response
`cat-picture = read(fd ...)`

Every HTTP library uses sockets under the hood



`$ curl awesome.com`

`Python: requests.get("yarl.us")`

oh, cool, I could write an HTTP library too if I wanted. Neat!
* SO MANY edge cases though!

AF-INET? What's that?

AF-INET means basically "internet socket": it lets you connect to other computers on the internet using their IP address.

The main alternative is AF-UNIX ("unix domain socket") for connecting to programs on the same computer.

3 Kinds of internet (AF-INET) sockets:

SOCK_STREAM = TCP
`curl uses this`

SOCK_DGRAM = UDP
`dig (DNS) uses this`

SOCK_RAW = just let me send IP packets.
`ping uses this`
I will implement my own protocol.

file buffering

15

!?!?
I printed some text but it didn't appear on the screen. why???

write

please write "I ❤️ cats"
process to file #1 (stdout)

okay!

Linux

time to learn about flushing!

3 kinds of buffering (defaults vary by library)

① None. This is the default for `stderr`.

② Line buffering. (write after newline). The default for `terminals`.

③ "full" buffering. (write in big chunks)
The default for `files` and `pipes`.

On Linux, you write to files & terminals with the system call

write

I/O libraries don't always call `write` when you print.

`printf("I ❤️ cats");`

wait for a newline
before actually writing

This is called **buffering** and it helps save on syscalls.

flushing

To force your I/O library to write everything it has in its buffer right now, call `flush`!

I'll call `write` right away!!
stdio

Python example:
`print("password:", flush=True)`

no seriously, actually write to that pipe program please

What's in a process?

12

PID	USER and GROUP	ENVIRONMENT VARIABLES	SIGNAL HANDLERS
process # 129 reporting for duty!	Who are you running as? julia!	like PATH! You can set them with: \$ env A=val ./program	I ignore SIGTERM! I shut down safely!
WORKING DIRECTORY	PARENT PID	COMMAND LINE ARGUMENTS	OPEN FILES
Relative paths (./blah) are relative to the working directory! chdir changes it.	PID 1 (init) is everyone's ancestor PID 147 PID 129	see them in /proc/PID/cmdline	Every open file has an offset. I've read 8000 bytes of that one
MEMORY	THREADS	CAPABILITIES	NAMESPACES
heap! stack! ≡ shared libraries! the program's binary! mmaped files!	sometimes one Sometimes LOTS	I have CAP_PTRACE Well I have CAP_SYS_ADMIN	I'm in the host network namespace I have my own namespace! container process

13

threads

Threads	Sharing memory (race conditions!)
Threads let a process do many different things at the same time	Threads in the same process share memory
process:	<p>I'll write some digits of π to 0x129420 in memory</p> <p>uh oh! That's where I was putting my prime numbers.</p>

Sharing memory (race conditions!)	RESULT: 24 ← Wrong. Should be 25! ←
memory ↑ I'm going to add 1 to that number! thread 1	<p>→ sharing data between threads is very easy. But it's also easier to make mistakes with threads.</p> <p>You weren't supposed to CHANGE that data!</p> <p>→ thread 1</p>