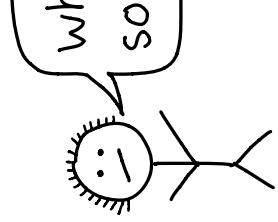




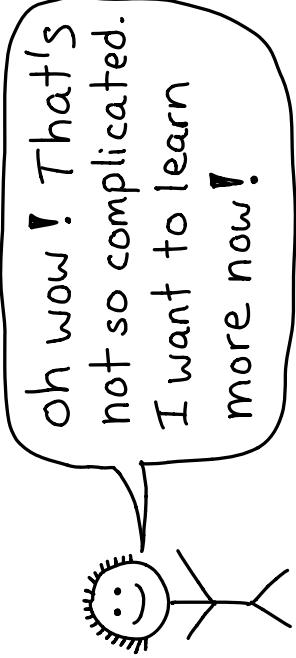
love this? ← find more awesome zines at wizazines.com →



what's a
socket?

you're in the
right place! This
zine has 19 comics
explaining important
Linux concepts.

... 5 minutes later ...



oh wow! That's
not so complicated.
I want to learn
more now!

by Julia Evans
<https://jvns.ca>
twitter.com/b0rk

Want to learn more?
I highly recommend
this book:

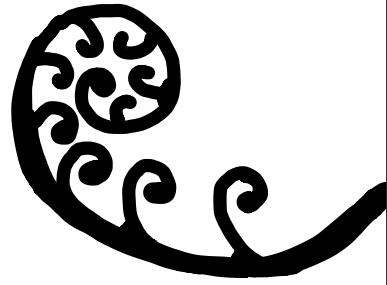
THE LINUX PROGRAMMING INTERFACE

MICHAEL KERRISK

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

\$ man printf

\$ man fopen

\$ man null

for /dev/null docs

④ devices

\$ man s|

⑤ file formats

\$ man sudoers

for /etc/sudoers

\$ man proc

files in /proc!

⑥ games

\$ man apt

not super useful.

⑦ miscellaneous

\$ man chroot

is my favourite

from that section

⑧ sysadmin programs

\$ man apt

♥ 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

→ **read** → **write** → **execute**

r = can read
w = can write
x = can execute

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

r → **rw-** → **r--** → **bork staff**

File permissions are 12 bits

setuid ↗ **setgid** ↗
000 ↗ **user** ↗ **group** ↗ **all**
sticky ↗ **110** ↗ **110** ↗ **100**
rw- ↗ **x** ↗ **rw-** ↗ **rw-**

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

`chmod 644 file.txt`
means change the permissions to:
r → **rw-** → **r--** → **rw-**
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! why?!

it's a long story ↗ unix ↗

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 of 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 ↗ dynamic linker ↗ I always mmap, so that file is probably 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	0x1a0000 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.

an amazing directory: /proc 5

Every process on Linux has a PID (process ID) like 42.

In **/proc/42**, there's a lot of VERY USEFUL information about process 42.

/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! ♡

Some pages are marked as either

- ★ read only
 - ★ not resident in memory
- your program stops running

→ Linux kernel code to handle the page fault runs

a page that's marked "not resident in memory," it triggers a **page fault!**

Linux running

how swap works

- ① run out of RAM
- ② Linux saves some RAM data to disk
- ③ mark those pages as "not resident in memory"
- ④ when a program tries to access the memory, there's a **Page Fault!**
- ⑤ Linux → time to move some data back to RAM!
- ⑥ if this happens a lot, your program gets **VERY SLOW** :-(I'm always waiting for data to be moved in & out of RAM

/proc/PID/cmdline
command line arguments the process was started with

/proc/PID/environ
all of the process's environment variables

/proc/PID/exe
symlink to the process's binary
magic: works even if the binary has been deleted!

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

/proc/PID/stack
The kernel's current stack for the process. Useful if it's stuck in a system call.
man proc
and `:more:`
Look at

for more information!

System calls

6

The Linux kernel has code to do a lot of things	your program <u>doesn't</u> know how to do those things	programs ask Linux to do work for them using system calls
read from a hard drive	make network connections	: TCP? dude I have no idea how that works
create new process	kill process	NO, I do not know how the ext4 filesystem is implemented. I just want to read some files!
change file permissions	Keyboard drivers	
every program uses system calls	I use the 'open' syscall to open files	So what's actually going on when you change a file's permissions is:
	me too!	run syscall #90 program with these arguments
	me three!	ok! Linux
	Python program	

copy on write

19

On Linux, you start new processes using the fork() or clone() system call.	calling fork creates a child process that's a copy of the caller	Linux does this by giving both the processes identical page tables.
	parent	RAM
	child	same → RAM
so Linux lets them share physical RAM and only copies the memory when one of them tries to write	I'd like to change that memory	but it marks every page as read only .
	Okay! I'll make you your own copy!	Linux

Shared libraries

18

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

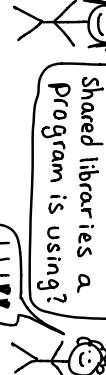
`openssl`
(for SSL!)

`sqlite`
(embedded db!)

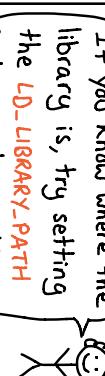
`libpcre`
(regular
expressions!)

`zlib`
(gzip!)

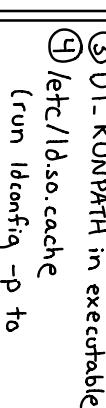
`libstdc++`
(C++ standard
library!)

 "How can I tell what shared libraries a program is using?"

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

 "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!
 dynamic linker

 Where the dynamic linker looks
 ① `DT_RPATH` in executable
 ② `LD_LIBRARY_PATH`
 ③ `DT_RUNPATH` in executable
 ④ `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"

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

`your code` `z lib` `sqlite`

`your code` `z lib` `sqlite`

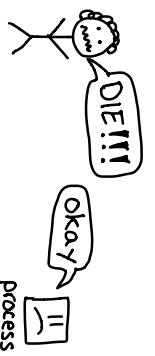
`your code` `z lib` `sqlite`

and programs like this
 ② Use separate shared libraries
`your code` `z lib` `sqlite` ← all different files

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

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!
 !! 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) (okay! I'll clean up and then exit!) oo process exceptions: SIGSTOP & SIGKILL can't be ignored got → !!

file descriptors

8

Unix systems use integers to track open files

`Open foo.txt` 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 → OS → fd 5!

write GET / HTTP/1.1 to fd #5 done!

`lsof` (list open files) will show you a process's open files

\$ ls -P 4242 ← PID we're interested in

FD NAME
0 /dev/pits/thy1
1 /dev/pits/thy1
2 pipe:29174 ← LOTS MORE (eventfd, inotify, signalfd, epoll, etc.)

not EVERYTHING on Unix is a file, but lots of things are

file descriptors can refer to:

- files on disk
- pipes
- sockets (network connections)
- terminals (like xterm)
- devices (your speaker! /dev/null!)
- LOTS MORE (eventfd, inotify, signalfd, epoll, etc.)

Let's see how some simple Python code works under the hood:

Python:

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

Behind the scenes:

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

(almost) every process has 3 standard FDs:

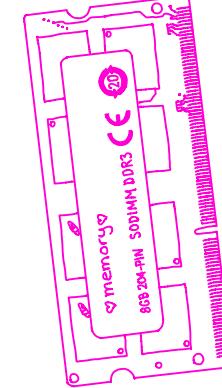
- stdin → 0
- stdout → 1
- stderr → 2

"read from stdin" means "read from the file descriptor 0", could be a pipe or file or terminal

virtual memory

17

your computer has physical memory



physical memory has addresses, like 0 - 8GB but when your program references an address like 0x5c69a2a2, that's not a physical memory address! It's a **virtual** address.

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

or I'll look that up in the page table and then access the right memory management unit, hardware

PID	Virtual Addr	Physical Addr
1971	0x20000	0x192000
2310	0x20000	0x228000
2310	0x21000	0x9788000

every program has its own virtual address space

program 1: 0x124520 → "puppies"

program 2: 0x129520 → "bananas"

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

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

Linux

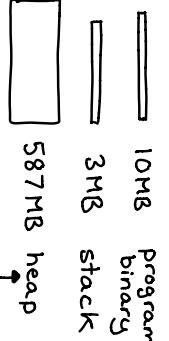
thanks, I'll use that now!

MMU

memory allocation

16

Your program has memory

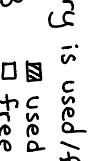


Your memory allocator's interface

`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 2 things.

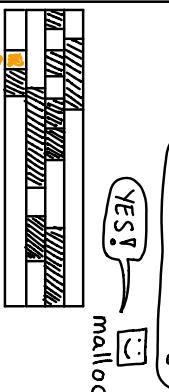
THING 1: Keep track of what memory is used/free.



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

Your code

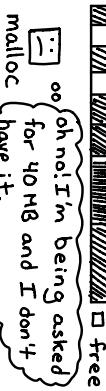
can I have 512 bytes of memory?



`malloc`

→ use a different `malloc` library like `jemalloc` or `tcmalloc` (easy!) → implement your own `malloc` (harder)

THING 2: Ask the OS for more memory!



`malloc`

can I have 60MB more?

here you go!

OS

`malloc` isn't magic! It's just a function!

you can always:

→ `use a different malloc library like jemalloc or tcmalloc (easy!)`
→ `implement your own malloc (harder)`

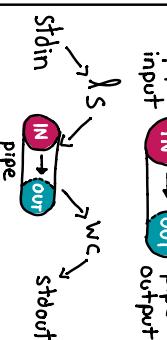
pipes

9

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

\$ ls | wc -l

53
53 files!



When `ls` does `write(IN, "hi")`, `wc` can read it!

`wc` can read it!
read(OUT)
→ "hi"

Pipes are one way: →
You can't write to OUT.

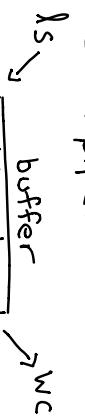
named pipes

\$ mkfifo my-pipe

This lets 2 unrelated processes communicate through a pipe!

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

Linux creates a buffer for each pipe.



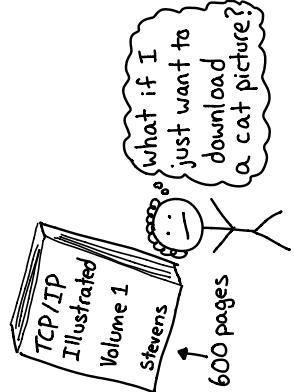
If data gets written to the pipe faster than it's read, the buffer will fill up. IN → OUT

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

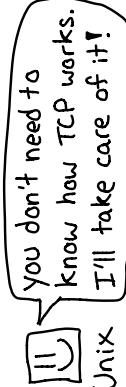
Sockets

10

networking protocols are complicated



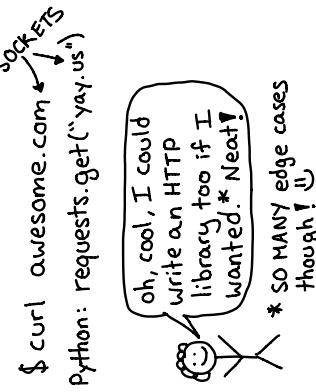
Unix systems have an API called the "socket API" that makes it easier to make network connections



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, (struct sockaddr_in){ .sin_family = AF_INET, .sin_port = htons(80) })`
- ③ Make a request
`write(fd, "GET /cat.png HTTP/1.1\r\n\r\n")`
- ④ Read the response
`cat-picture = read(fd, ...)`

Every HTTP library uses sockets under the hood



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.

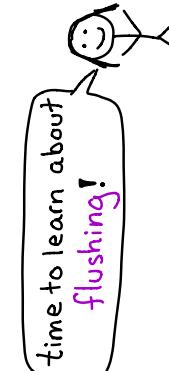
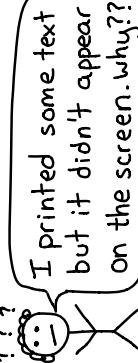
The main alternative is AF-UNIX ("unix domain socket") for connecting to programs on the same computer.

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 I will implement this my own protocol.

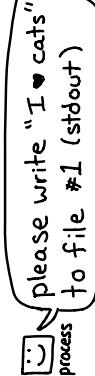
file buffering

15



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

write ❤



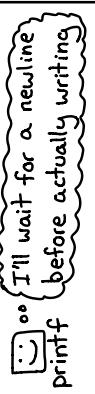
flushing

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

- 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**.

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

`printf("I ❤ cats");`



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

when it's useful to flush
→ when writing an interactive prompt!

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

→ when you're writing to a pipe / socket
no seriously, actually write to that pipe please

what's in a process?

12

process?

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

13

threads

Threads let a process do many different things at the same time	and they share code	why use threads instead of starting a new process? → a thread takes less time to create.
process:	calculate-pi find-big-prime-number	sharing data between threads is very easy. But it's also easier to make mistakes with threads.
Thread 1	I'll write some digits of π to 0x129420 in memory	You weren't supposed to CHANGE that data!
Thread 2	uh oh! that's where I was putting my prime numbers.	min thread 1 CPU 1 π thread primes thread