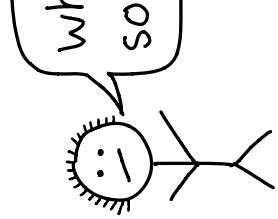


By Julia Evans

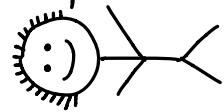
love this?
find more awesome zines at
wizardzines.com ←



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

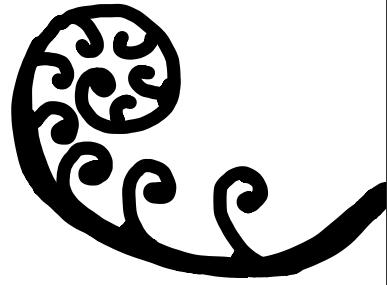
THE LINUX PROGRAMMING INTERFACE

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

\$ man printf

\$ man null

\$ man fopen

for /dev/null docs

④ devices

\$ man ls

not super useful.

⑤ file formats

\$ man sudoers

\$ man proc

is my favourite
from that section

⑥ games

\$ man apt

\$ man chroot

⑦ miscellaneous

explains concepts!

\$ man 7 pipe

\$ man 7 symlink

♥ Table of contents ♥

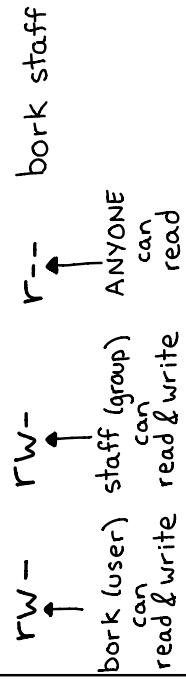
unix permissions...4	sockets.....10	virtual memory...17
/proc.....5	unix domain sockets.....11	shared libraries...18
system calls.....6	processes.....12	copy on write....19
signals.....7	threads.....13	page faults.....20
file descriptors.8	floating point....14	mmap.....21
pipes.....9	file buffering.....15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

Unix permissions

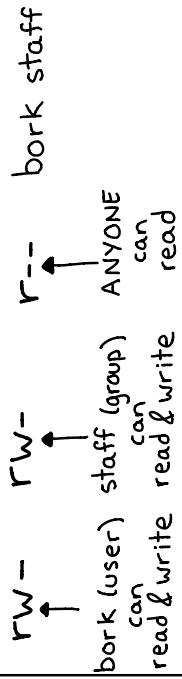
4

There are 3 things you can do to a file

→ read write execute



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



File permissions are 12 bits

setuid setgid → 000 user group all
sticky → 110 110 100
rw x → 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

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.



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!

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!

↪ I need to use libc.so.6 C standard library

↪ you too eh? no problem. 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	0x1a000 readonly

"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! ♡

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

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!

how swap works

- ① run out of RAM → RAM full
 - ② Linux saves some RAM data to disk
 - ③ mark those pages as "not resident in memory" in the page table
 - ④ when a program tries to access the memory, there's a !Page fault!
 - ⑤ :-(+ 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

what happens during a page fault?
→ the MMU sends an interrupt

your program stops running

Linux kernel code to handle the page fault runs

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

System calls

The Linux kernel has code to do a lot of things	your program <u>doesn't</u> know how to do those things	make network connections	TCP? dude I have no idea how that works	programs ask Linux to do work for them using <code>system calls</code>
read from a hard drive	NO, I do not know how the ext4 filesystem is implemented. I just want to read some files!	kill process	Please write to this file	program <switch to running kernel code>
create new process	So what's actually going on when you change a file's permissions is:	change file permissions	done! I wrote 1097 bytes!	Linux <program resumes>
every program uses system calls	run syscall #90 with these arguments	and every system call has a number (e.g. chmod is #90 on x86-64)	ok!	Linux
I use the 'open' syscall to open files	me too!	I use the 'open' syscall to open files	me three!	Python program
Java program	C program			

Copy on write

On Linux, you start new processes using the <code>fork()</code> or <code>clone()</code> system call.	calling fork creates a child process that's a copy of the caller	the cloned process has EXACTLY the same memory.	copying all that memory every time we fork would be slow and a waste of RAM	when a process tries to write to a shared memory address:
		→ same heap	often processes call exec right after fork, which means they don't use the parent process's memory basically at all!	① there's a <code>≈page fault</code> :
		→ same stack		② Linux makes a copy of the page & updates the page table
		→ same memory maps if the parent has 3GB of memory, the child will too.		③ the process continues, blissfully ignorant
			RAM	It's just like I have my 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 !!

ldd!!!

I got a "library not
found" error when running
your code

↳ all different
files

If you know where the
library is, try setting
the LD_LIBRARY_PATH
environment variable

↳ LD_LIBRARY_PATH
dynamic linker
tells me where to look!

There are 2 ways
to use any library:
① Link it into your binary

[your code] [zlib] [sqlite]
big binary with lots of things!

Programs like this
are called "statically linked"
and programs like this

[your code] [zlib] [sqlite]

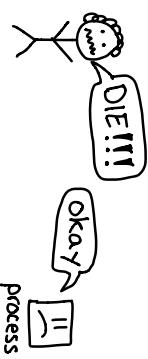
Where the dynamic
linker looks

① DT_RPATH in 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

Signals

7

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



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

SIGINT ctrl-c SIGTERM kill
SIGTERM kill -9 } various
levels of
"die"

SIGHUP kill - HUP
often interpreted as
"reload config", e.g. by nginx

Your program can set
custom handlers for
almost any signal

your child
that pipe
is closed
terminated
illegal
instruction
the timer you
set expired
Segmentation
fault

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

process

① ignore

② kill process

③ kill process AND
make core dump file

④ stop process

⑤ resume process

SURPRISE!
another signal!

Every signal has a default
action, which is one of:
① ignore
② kill process
③ kill process AND
make core dump file
④ stop process
⑤ resume process

SIGKILL (terminate)
SIGSTOP & SIGKILL
can't be ignored
got → [x]

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

ls -l (list open files) will show you a process's open files

\$ ls -l -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.)

3 /home/bark/awesome.txt
4 /tmp/
5

FD is for file descriptor

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.)

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

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

Python:
f = open("file.txt")
f.read_lines()

Behind the scenes:

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

(almost) every process has 3 standard FDs:

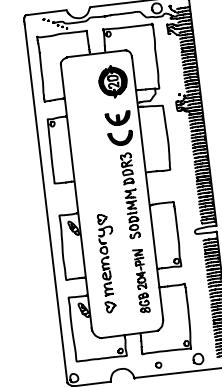
- stdin → 0
- stdout → 1
- stderr → 2

"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.

every program has its own virtual address space

0x129520 → "puppies"
program 1

0x129520 → "bananas"
program 2

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

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

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

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

Linux

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

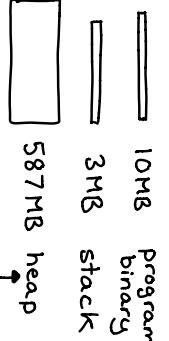
MMU

thanks, I'll use that now!

memory allocation

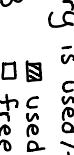
16

your program has memory



your memory allocator (malloc) is responsible for 2 things.

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



the heap is what your allocator manages

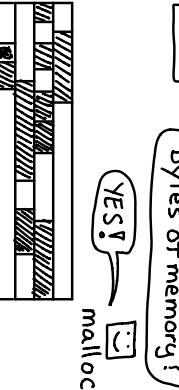
your memory allocator's interface

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.

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

Your code

can I have 512 bytes of memory?



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)

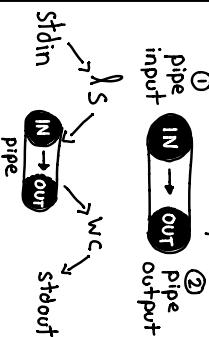
pipes

9

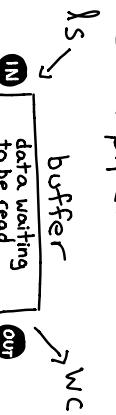
Sometimes you want to send the output of one process to the input of another

```
$ ls | wc -l
```

```
53
53 files!
```



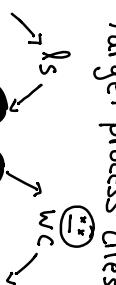
Linux creates a buffer for each pipe.



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

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

what if your target process dies?



named pipes

\$ mkfifo my-pipe

This lets 2 unrelated processes communicate through a pipe!

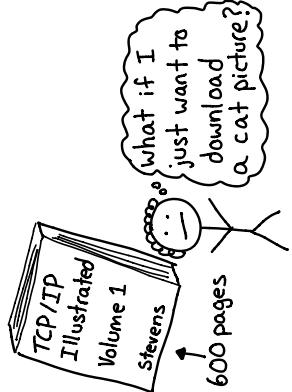
If wc dies, the pipe will close and ls will be sent SIGPIPE. By default, SIGPIPE terminates your process.

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

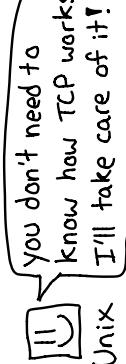
Sockets

10

networking protocols are complicated



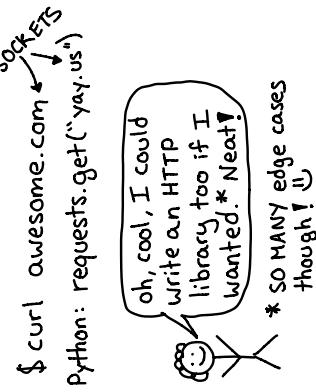
Unix systems have an API called the "socket API" that makes it easier to make network connections



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_addr.s_addr = 12.13.14.15, .sin_port = 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



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

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

time to learn about flushing!

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

write *

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

Okay! Linux

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

printf("I ❤ cats");

:)" I'll wait for a newline printf before actually writing

This is called buffering and it helps save on syscalls.

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.

flushing

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

no seriously, actually write to that pipe please

floating point

14

a double is 64 bits

Sign exponent fraction
 \downarrow \downarrow \downarrow
 100101 100101 100101 100101
 100101 100101 100101 100101

$$\pm 2^{E-1023} \times 1.\text{frac}$$

That means there are 2^{64} doubles.
The biggest one is about

$$2^{1023}$$

doubles get farther apart as they get bigger

between 2^n and 2^{n+1} there are always 2^{52} doubles, evenly spaced.

that means the next double after 2^{60} is $2^{60} + 2^{64} = \frac{2^{60}}{2^{52}}$

$$2^{52} + 0.2 = 2^{52}$$

\leftarrow (the next number after 2^{52} is $2^{52} + 1$)

$$1 + \frac{1}{2^{54}} = 1$$

\leftarrow (the next number after 1 is $1 + \frac{1}{2^{52}}$)

$$2^{2000} = \text{infinity}$$

\leftarrow infinity is a double

$$\text{infinity} - \text{infinity} = \text{nan}$$

\leftarrow nan = "not a number"

Javascript only has doubles (no integers!)

$$> 2^{**53}$$

9007199254740992

$$> 2^{**53+1}$$

9007199254740992

$$\begin{matrix} & \nwarrow \\ \text{same number!} & \text{uh oh!} \end{matrix}$$

doubles are scary and their arithmetic is weird!

they're very logical! just understand how they work and don't use integers over 2^{53} in Javascript

unix domain sockets !!

unix domain sockets are files.

```
$ file mysock.sock
```

socket

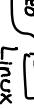
the file's permissions determine who can send data to the socket.

advantage 1

lets you use file permissions to restrict access to HTTP/ database services!

chmod 600 secret.sock

This is why Docker uses a unix domain socket.

 run evil container permission denied 

advantage 2

UDP sockets aren't always reliable (even on the same computer). unix domain datagram sockets are reliable! And they won't reorder packets!

 I can send data from sketchy.com  video decoder  Sandboxed process

advantage 3

You can send a file descriptor over a unix domain socket. Useful when handling untrusted input files!

 here's a file I downloaded from sketchy.com  video decoder  Sandboxed process

what's in a process?

12

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

13

threads	Threads let a process do many different things at the same time	and they share code
	processes share memory	calculate-pi find-big-prime-number
	I'll write some digits of π to 0x129420 in memory	but each thread has its own stack and they can be run by different CPUs at the same time
	uh oh! that's where I was putting my prime numbers.	cpu 1 CPU 2 π thread Primes thread
Sharing memory can cause problems (race conditions!)	Thread 1: I'm calculating ten million digits of π! so fun! Thread 2: I'm finding a REALLY BIG prime number!	Why use threads instead of starting a new process? → a thread takes less time to create.
	memory ↑ 2 3 I'm going to add 1 to that number!	sharing data between threads is very easy. But it's also easier to make mistakes with threads.
	thread 1: I'm going to add 1 to that number!	thread 2 ↓ memory ↑ 2 3 I'm going to add 1 to that number!
		You weren't supposed to CHANGE that data!

RESULT: 24 ← WRONG. Should be 25!