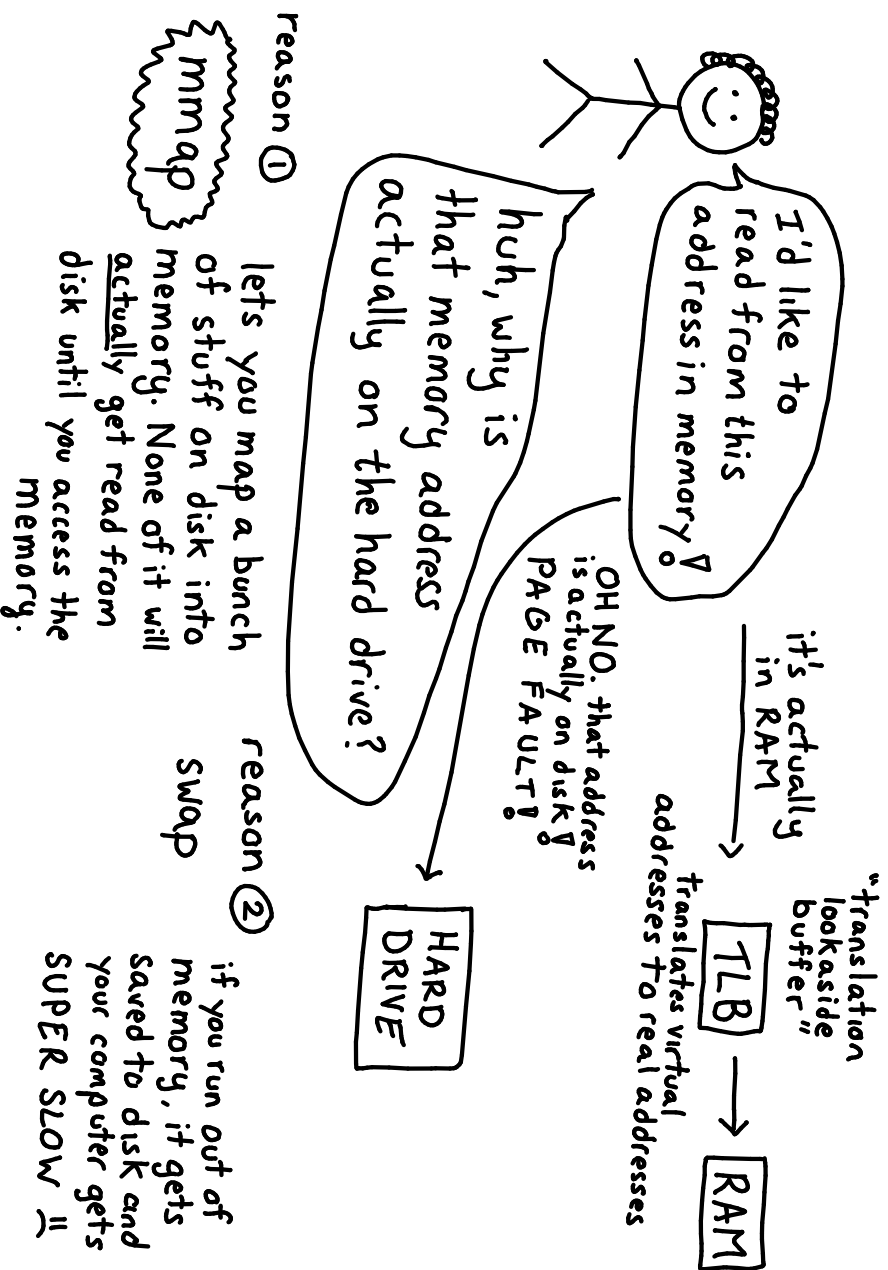


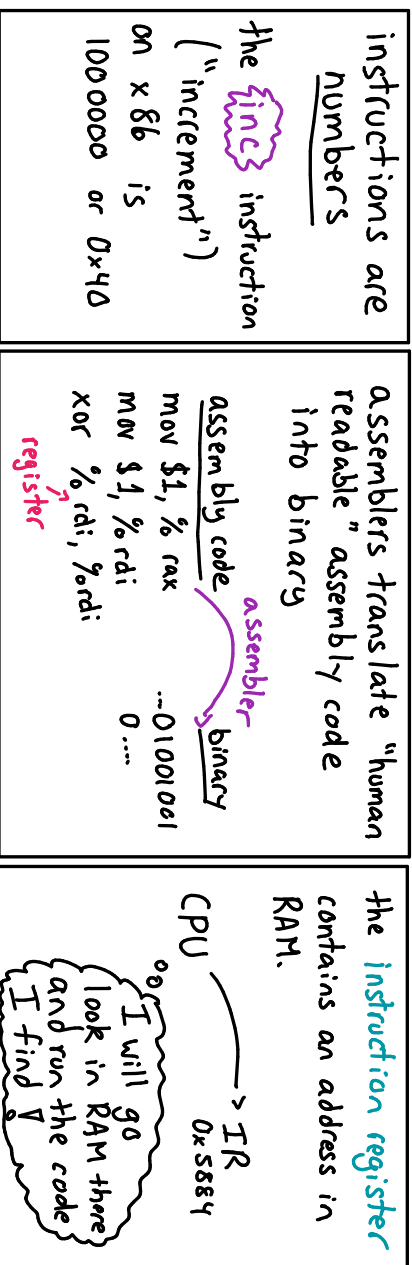
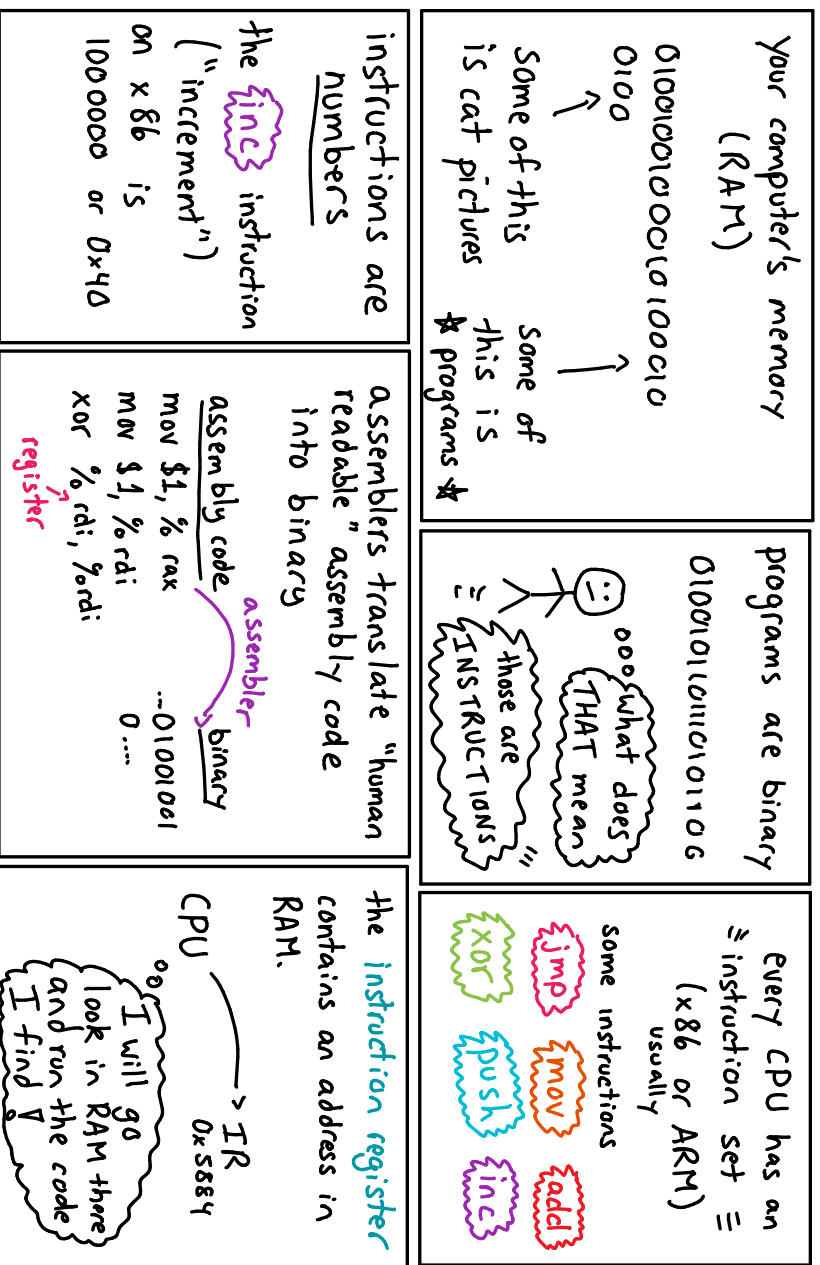
≡ VIRTUAL MEMORY ≡



assembly

SUIA EVANS
@b0rk

We hear computers "think in binary". But what does that MEAN??



"C" stands for linearizable

the CAP theorem

Sulia Evans
@bork

from Martin Kleppmann's "A critique of the CAP theorem"

in distributed systems,
network partitions happen
???

hello?

computer

someone unplugged a cable!

garbage collection!
too much network traffic!

if you want to be
consistent you can't
always be available

panda

elephant

You're gonna have
to wait for an
answer

"CP systems"

consul

etcd

zookeeper

chubby

when they reply, you can
believe them, but they
don't always give you
answers

"AP systems" available +
partition tolerant
this doesn't mean very much.

always return "lol"

very carefully
considered weaker
consistency model

You can call both of
these "AP"

CAP is a
very simple theorem

I read the
whole proof!
It took
10 minutes +
there's no
math

CAP won't help
you reason about
most systems

I have a
replicated database
what can you
tell me?

nothing!

CAP

drawings.jvns.ca

User space vs. kernel space

Sulia Evans
@bork

the Linux kernel has
millions of lines of code

- ★ read + write files
- ★ decide which programs
get to use the CPU
- ★ make the keyboard
work

When Linux kernel
code runs, that's
called

kernel space

When your program
runs, that's

user space

time to
write a file

that's MY
JOB!

YOUR PROGRAM

KERNEL

time for a
Context switch
to kernel space

Your program switches
back and forth
str = "my string"

x = x + 2

file.write(str) ← ★ switch to
kernel space

y = x + 4

str = str * y ← ★ and we're
back to
user space!

timing your process

\$ time find /home

0.15 user 0.73 system

time spent in
your process

time spent by
the kernel doing
work for your
process

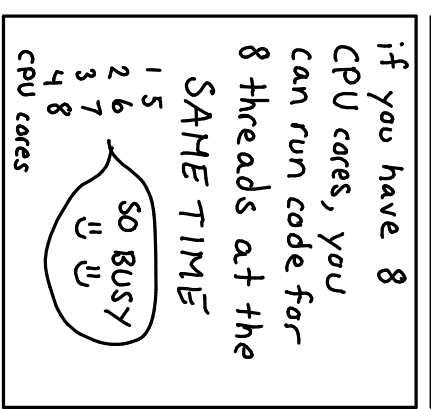
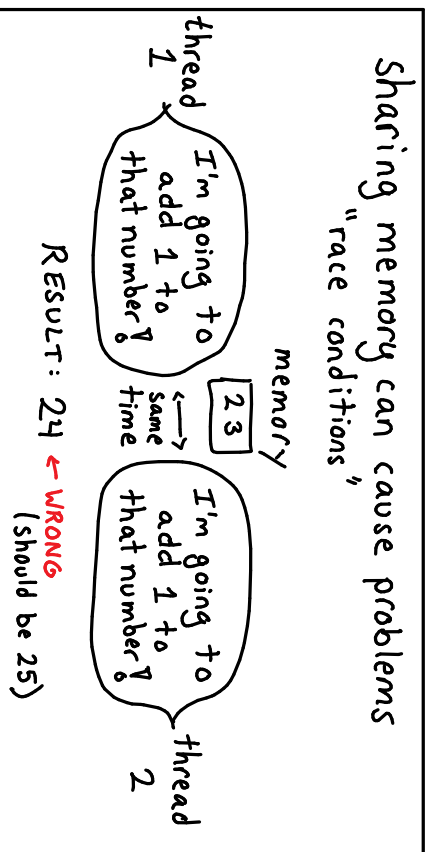
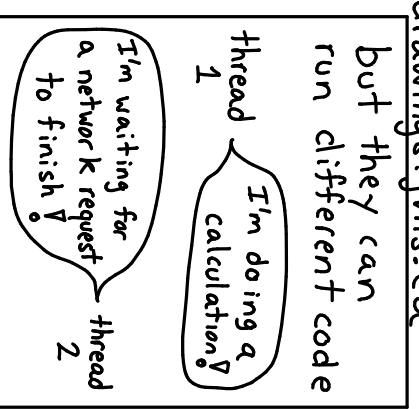
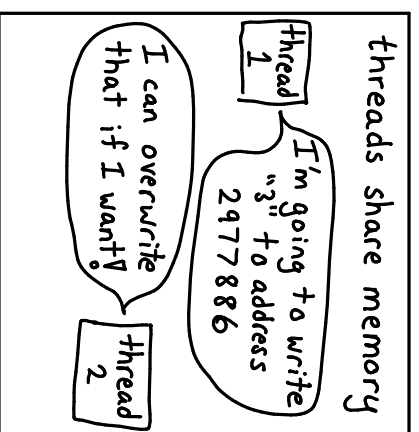
What's a thread?

JULIA EVANS
@b0rk

drawings.jvns.ca

a process can have lots of threads

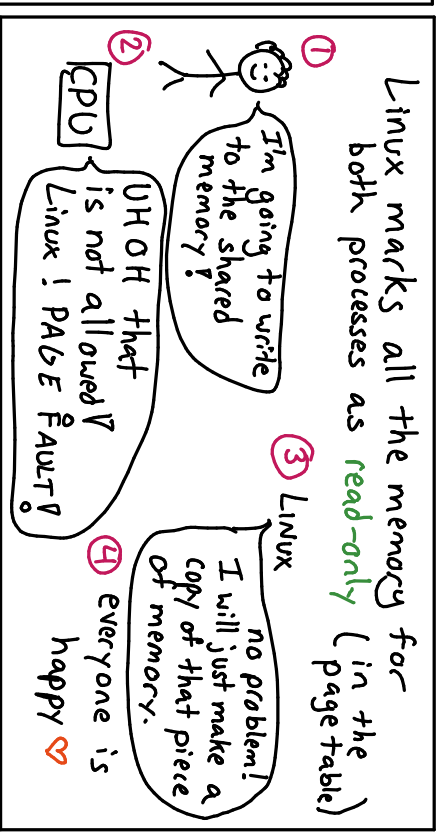
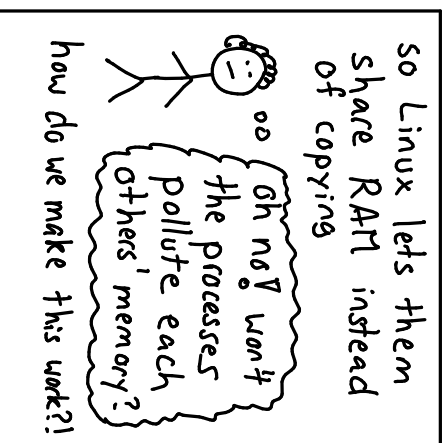
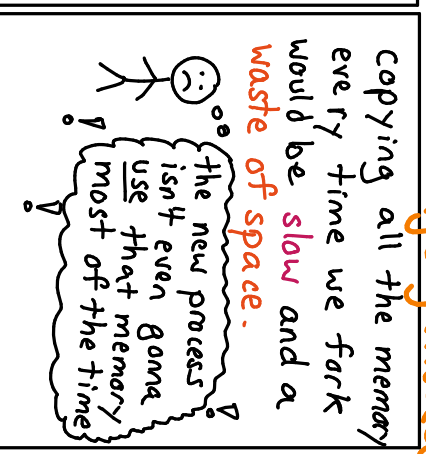
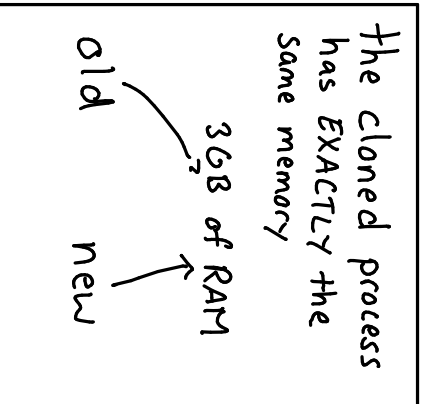
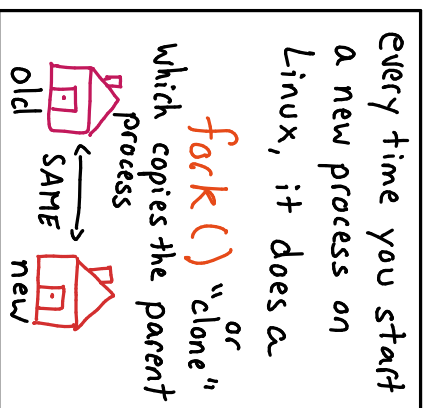
process id	thread id
1888	1888
1888	1892
1888	1893
1888	2007



JULIA EVANS
@b0rk

copy on write

drawings.jvns.ca



directories + symlinks

@b0rk
Julia Evans
drawings.jvns.ca

What's a directory?

filename	inode	number
awesome.jpg	279932	
blah.txt	13227	
cumberbatch	238333	

What's a symlink?

it's just a file with the name of another file in it

\$ **readlink** my-cool-link

/home/julia/long-complicated-file-name

I made a directory with 2,000,000 files

It's so SLOW

Listing
Your directory
is gonna be
REAL SLOW

(a few seconds at least)

OLD

on ext 2 even opening files in big directories is slow

that's right! ext 2 directories have no index so you have to SEARCH THE WHOLE THING

ext 2 is OLD though. ext 3 is OK.

more at
drawings.jvns.ca

★ the stack ★

(in a C program)

JULIA EVANS
@b0rk

your program has

→ local variables

```
int x = 2;
```

→ a function to return to

```
void parent() {
```

```
    do_thing();
```

→ function arguments

```
make_cat(name, fluffiness)
```

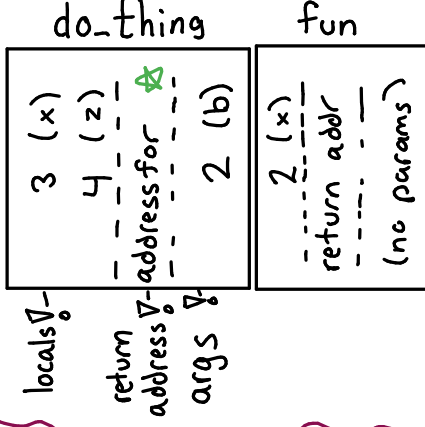
these all live in a part of memory called

the stack

example program

```
int fun() {  
    int x = 2;  
    do_thing(2);  
    int y = 4;   
    void do_thing(b) {  
        int x = b+1;  
        int z = 4;  
        return;  
    }  
}
```

the stack at



there's a limit to how big your stack can get! Exceed it and you get a

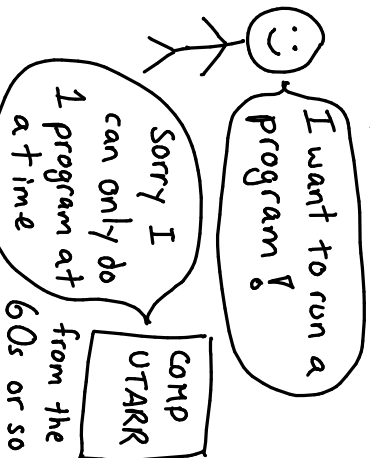
STACK OVERFLOW

drawings.jvns.ca

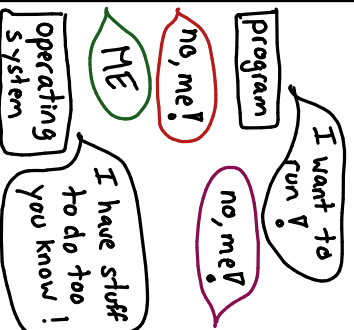
CPU scheduling

JULIA EVANS
@b0rk

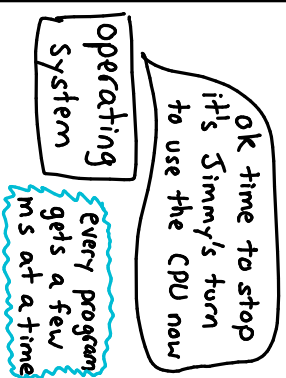
once upon a time...



your computer today



every CPU core can only run 1 program at a time



steps when we switch the running process

"context switch"

- save:
 - registers
 - stack pointer
 - which CPU instruction to start at next time
- set up memory for new process
- load new registers and stuff

all this takes time
(2 microseconds?).
It's ok to do but
you don't want to
be switching processes
constantly

you don't use the CPU when you're waiting

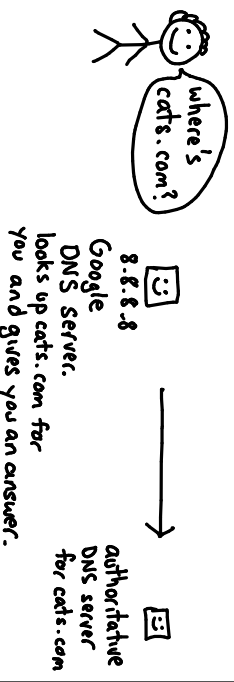


how does DNS work?

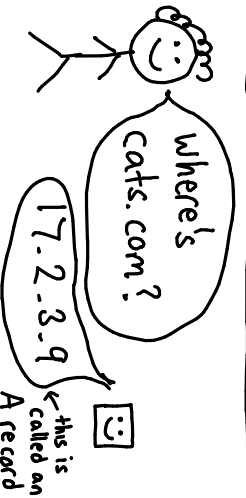
JULIA
EVANS
@b0rk

more of these at drawings.jvns.ca

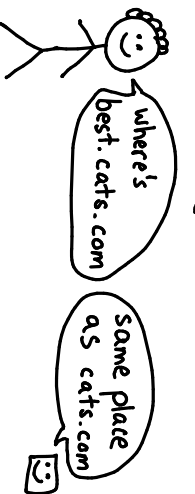
most DNS queries get cached



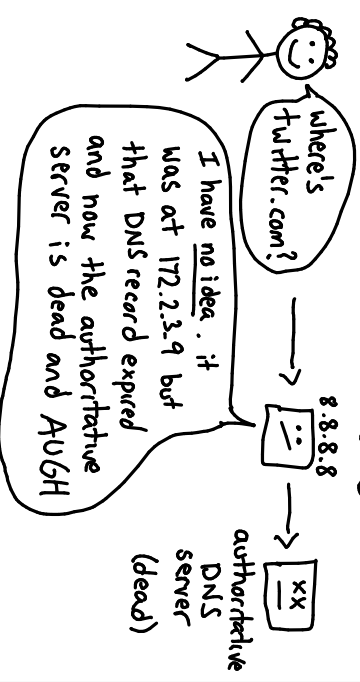
DNS servers translate names to IP addresses



sometimes they tell you it's an alias (CNAME record)



when an important DNS server dies



floating point

SULIA EVANS
@bork
more at: drawings.jvns.ca

a double is 64 bits.

that means there are 2^{64} different doubles

going up to

1.8×10^{308}

some double arithmetic

$2^{52} + 0.2 = 2^{52}$ (the next number after 2^{52} is $2^{52} + 1$)

$1 + \frac{1}{2^{53}} = 1$ (the next number after 1 is $1 + \frac{1}{2^{52}}$)

$3 \times 10^{100} = \text{infinity}$ \leftarrow infinity is a double
infinity - infinity = nan (not a number)

there are 2^{52} numbers between 1 and 2

$1 + \frac{1}{2^{52}}, 1 + \frac{2}{2^{52}}, \dots$

2^{51} numbers between 2^{54} and 2^{55}
between 4 and 8
etcetera.

JavaScript only has doubles (and Lua!)

that means after 2^{53} you don't have every integer!

printing doubles is nontrivial

the shortest version of 25.64853898042e8 is 2.564854e9
calculating the shortest representation takes time!

SULIA EVANS
@bork

asking good questions

find a good time

hey can I ask you about database performance for 20 minutes?

yeah! can we do it after lunch?

yeah!

state what you know

so, I know when the database gets a lot of writes, the hard drive can't keep up.

that's right! I don't think that was our problem though, look at this...

ask factual questions

does this database take out a lock when it does writes?

yes! here are the docs you should read if you want to know more! They're good.

choose who to ask

probably a better choice, has a good shot at answering your questions + way more time

the database creator

your coworker with a bit more experience than you

do some research

so I found out that creating database indexes takes time and I have questions about how that affects performance...

great

profit

now I know a lot more

I really helped! That was a great use of time

pipes

Sulia Evans
@bork

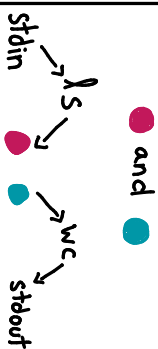
drawings.jms.ca

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

\$ ls | wc -l

53
← 53 files

a pipe is a pair of 2 magical file descriptors



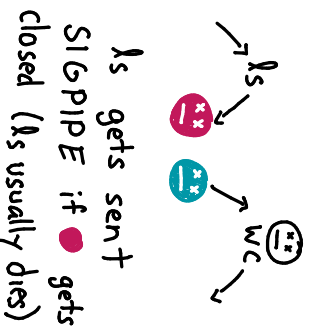
when **ls** does `writeln(●, "hi")`
wc can read it!
`read(●)`
→ "hi"

pipe buffers

Is I'm gonna write a bajillion bytes to ●

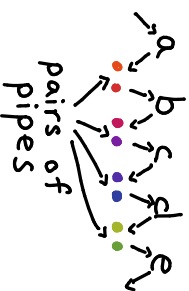
Uh no if my buffer is full you have to wait

what if your target process dies?



you can pipe SO MANY things together

\$ a | b | c | d | e



IPv6

Sulia Evans
@bork

drawings.jms.ca

hello can I have an IP address like 172.99.7.3?

you're OUT??

no we're out
yeah there were only 2³² (4 billion) to start with and there are SO MANY cell phones

What will we DO??

IPv4 addresses are 32 bits

you can have this IPv6 address though! (128 bits)

2001:04b8:85a3:0000:0000:8a2e:0370:27a9

that's cool! I totally understand IPv6 because this is 2016

Windows 2000 had IPv6 support. operating systems: READY

hello can I use this website

I am only set up for IPv4

Server

IPv6 user

Sometimes people put translators in the middle to turn IPv6 packets into IPv4 packets

adoption
it's happening

Google says 30% of American traffic they see is using IPv6

people were putting it off but we're REALLY RUNNING OUT of IPv4 addresses so now they have no choice

the "OSI model" for networking

SULIA EVANS
@bork

I don't always find it useful but it's good to know what "layer 4" means

what does "this is an L4 proxy" mean?

If a load balancer is labelled "L7", it usually means it looks at the Host: header inside your HTTP packets.

LAYERS

- 1: electrical engineering stuff, wires, frequencies, wifi
- 2: Ethernet protocol + others
- 3: IP (IP addresses)
- 4: TCP + UDP (ports)
- 5+6: nobody ever talks about these
- 7: HTTP and friends

layer 3
networking
tool

↑
ignores layer 4 and above

I only know about IP addresses! I don't even know what a port is let alone what the packet says

unix permissions

SULIA EVANS
@bork

3 kinds of things you can do to a file
↓ read ↓ write ↓ execute

\$ ls -l /bin/ping
rw[↓]x-r-x root root
setuid flag

This means ping always runs as root (who owns it), no matter who started ping

\$ ls -l awesome.png
rw- rw- r-- bork staff
↑ ↑
bork can do this (user) staff can do this (group)
ANYONE can do this

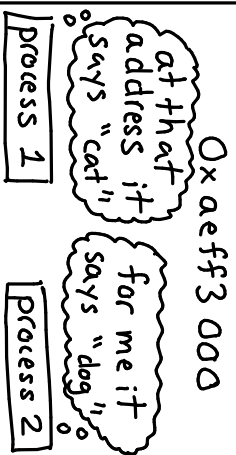
what's this 755 business?
7 means rwx
6 → rw-
5 → r-x
4 → r--
it's binary?
5 → 101 → r-x
755 means
rwx r-x r-x

more weird permissions things
setgid
sticky bit
but I ran out of space

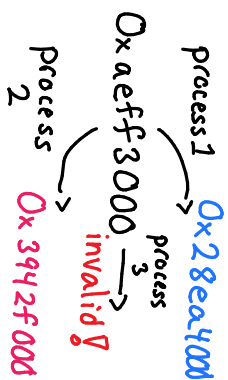
Sulia Evans
@bork

Page table (in 32 bit memory)

every process has its own memory space



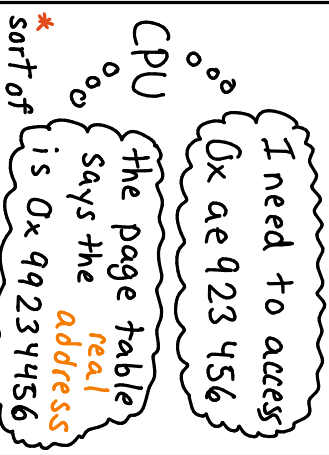
each address maps to a 'real' address in physical RAM



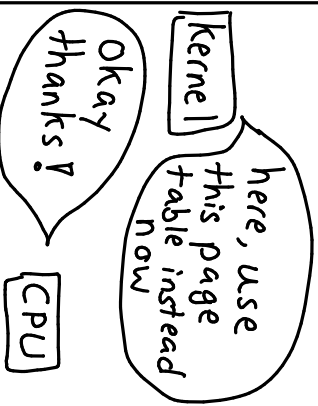
Processes have a "page table" in RAM that stores all their mappings

$0x12345000 \rightarrow 0xae925...$
 $0x23f49000 \rightarrow 0x12345...$
the mappings are usually 4kB blocks (4kB is the normal size of a "page")

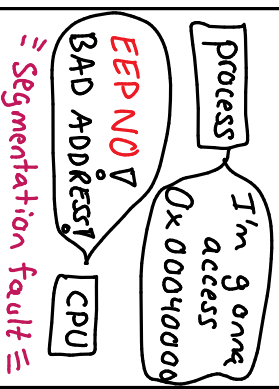
every *memory access uses the page table



when you switch processes...



some pages don't map to a physical RAM address

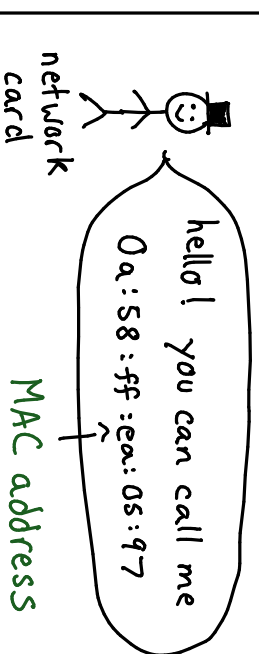


Sulia Evans
@bork

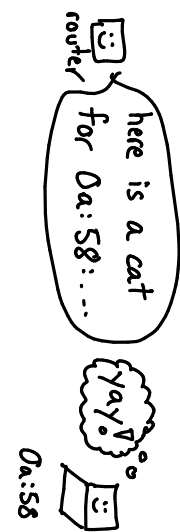
what's a MAC address?

more at: [drawings-jvns.ca!](https://drawings-jvns.ca/)

every computer on the internet has a network card



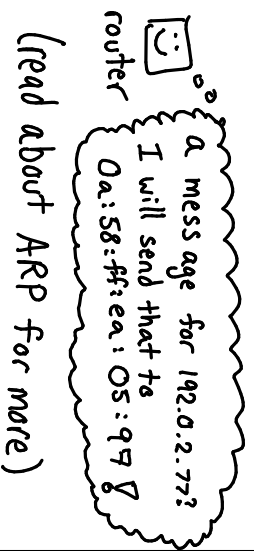
when you make HTTP requests with Ethernet/WiFi, every packet gets sent to a MAC address



wait, how do I know someone else on the same network isn't reading all my packets?

You don't! that's one reason we use HTTPS + secure WiFi networks

your router has a table that maps IP addresses to MAC addresses



memory allocation

Julia Evans
@b0rk

at any given time
your program has
a fixed amount of
memory

587 MB used free

this was used but then
got freed

and it can ask
the OS for more
memory

! now I have
1.8 GB of
memory! yay!

google
chrome

your allocator tries to
fill in unused pieces when
you ask for memory

can I have
512 bytes of
memory?

YES

malloc

↑ your new memory

you can invent your
own strategy to allocate
memory

glibc malloc's algorithm
is dumb I'm going to
do my own thing

this is sort of normal to do
if you care a LOT about performance

especially if you understand
how memory
allocation
works

alternatives to libc malloc

jemalloc
Facebook

tcmalloc
Google

SULIA EVANS
@b0rk

more at
drawings.jvns.ca

anatomy of a packet

When you get
a webpage like
Facebook, it
comes into your
computer in man-
small packets

Let's see what
those look like!

Packets are split into
a few sections
(or "headers")

← "physical layer".
this gets changed
constantly as your
packet moves between
computers.

← in charge of
getting your packet
to the right server
(like an address on
an envelope)

← in charge of
preventing data
corruption and helping
you retry lost packets

video streaming uses
UDP instead. UDP
does not try to be
reliable.

← the actual
data you're
trying to send!

ethernet/wifi:

82:53:ac:99:2f:33 ← MAC
address

IP ("internet protocol")

FROM: 172.96.2.3 TO: 123.9.2.32

TCP (or UDP)

sequence number: 877392 ← counts
bytes
checksum: 8447 ← detect corrupted
data sent so
far

from: port 9979 to: port 80

HTTP (or whatever)

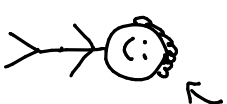
GET / HTTP/1.1
Host: google.com
Accept-Language: en-US

networking concepts

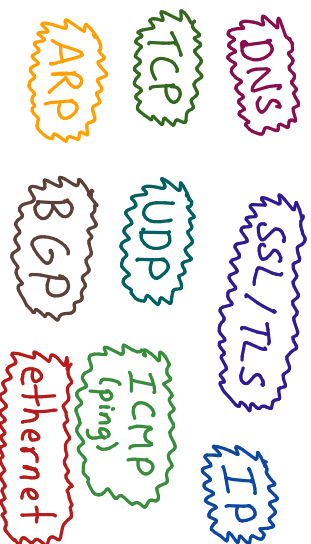
Sulia Evans
@bbark

hey I want to understand all the networking stuff that happens when I go to google.com!

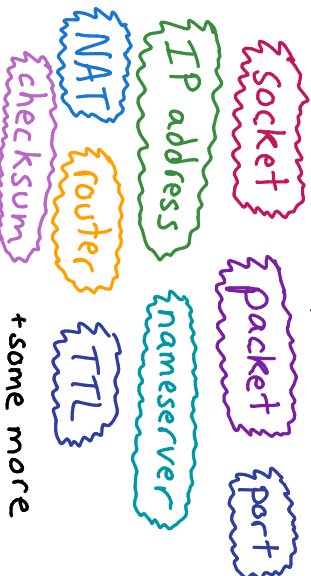
YES that is awesome. there are a lot of concepts but you can totally learn them all!



protocols



other concepts



it's a lot to learn but it's totally possible to learn how it all fits together to get you pictures of cats ☺

(knows many networking concepts now)

man pages = awesome

Sulia Evans
@bbark

I found out I can get documentation for programs (like grep) with **man grep**!

but that's not all!! lots of other things have man pages too!



man pages are split up into 8 sections

- ① ② ③ ④ ⑤ ⑥ ⑦ ⑧

/usr/share/man/man5 has section 5 on my machine.

(sometimes quality may vary ☹)

① programs

\$man grep
\$man ls

② system calls

\$man sendfile

③ C functions

\$man 3 printf
\$man fopen

④ devices

\$man null
for /dev/null docs

⑤ file formats

\$man sudoers
for /etc/sudoers
GREAT → \$man proc

⑥ games

(not very useful)
man sl is good if you have sl though

⑦ miscellaneous

\$man 7 pipe
\$man 7 symlink
(these are cool!)

⑧ sysadmin programs

\$man apt
\$man chroot

mesos

JULIA EVANS
@bork

mesos manages resources

master

agents

We have 200 CPUs + 800 GB of RAM. What should we do?

agents run "tasks"

running on agent #999 needs 2 GB of RAM + 3 CPUs

program ← state: running

the Mesos master keeps track of EVERY running task

!!!

dude there are THOUSANDS of these things. I got it though.

frameworks ask the Mesos master to run tasks

there are LOTS.

Chronos (cron-like jobs)
Marathon (HTTP services)
Senkins
Spark
Hadoop
ElasticSearch
Cassandra

you can split your Mesos cluster between several frameworks

half for Hadoop, half for web services!

Mesos doesn't know much about tasks

task

idk what it's doing

that's a HTTP service running on port 9923

mesos

Marathon

mutexes

JULIA EVANS
@bork

drawings.jvns.ca

Sometimes you're running code on 2 CPUs at the same time

x=2 CPU 1

x=3 CPU 2

Sometimes 2 threads want to change the same thing

array

write "a" write "u"

"hallo!" "hallo!"

CHAOS!

program 1

program 2

a mutex keeps track of whether something is in use

program 1's turn

Write "a"!

help not my turn.

program 2

when you're done, you tell the mutex it's available

mutex

I'm done!

Yay!

program 1

program 2

there's lots more but we're outta space

semaphores

mutex

Compare and swap

atomic instructions