Lecture 1

Red Star 3.0 (from fedora)
- automatic watermarking (steganography)
- tamper-resistant OS modification

http://web.cs.ucla.edu/classes/winter16/cs111

Course organization and grading

17 lectures
1/9   1 midterm (during lecture)     100 minutes
2/9   1 final exam     180 minutes
1/8   4 labs (teams of 2) ½ each
      Shell, kernel hacking, file system, networking
      miscellaneous
2/15    individual minilabs 1/15 each
      Scheduling, virtual memory    with partner
1/12   1 design problem (lab extension), requires written report
1/15   1 2-3 page paper on an operating system topic
1/20    scribe notes (groups of up to 4), webpage for lecture
      Due week after lecture

Open book/notes/assignments on the exams ☺

What's a System?
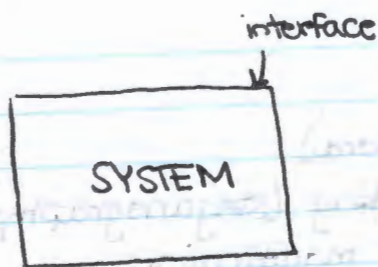   OED original (1928)
     I. An organized or connected group of objects
     II. A set of principles, etc., a scheme, method
   from Greek σύστημα' - organized whole,
     government, constitution, a body of people or animals,
     musical interval
  roots:
   - set up with"

interface

SYSTEM

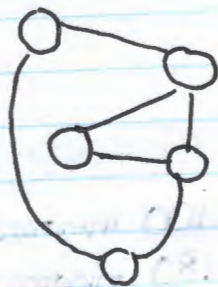ENVIRONMENT

Book's definition:
- A set of interconnected components that has a specified behavior observed at the interface of its environment.

Example of subsystems, not seen from environment.

Operating System - American Heritage Dict. (2000)

Software designed to control hardware of a specific data processing system in order to allow users and application programs to make use of it.

Encarta (2007)
master control program in a computer

Wikipedia v698 216 816 (2016 - 01 - 04)
collection of smaller programs and software that is used to control and operate the computer system
....

Some major missing issues...
• resource management
• reliability + error handling
• security

circa 2008 hardware

$ ls -l big

$\sim 10^{19}$

-rw-rw-r-- 1 eggert faculty 9223372036854771500
                -Oct 6  11:31 big ← ZFS
                                    Zeta File System

$ grep x big                                    $10^{21}$ bytes/second
$ time grep x big
    real  0m0.009s  $10^{-2}$ s              $10^{22}$ bits/second

http://what-if.xkcd.com/31/ 167
            Internet bandwidth. ~~167~~ Tb/s ~ $10^{14}$ bits/second

This means it's possible to ~~physically~~ move data faster
than it is to send it through the Internet, given a sufficient
size

$1.2 million/gallon → all freight trucks in the U.S.
        ~ 0.52 ~~bytes~~/second
               bits

ZFS has intensional files - the data is not stored on the
disk, rather, it is a description or a program that represents
the data.                              = It also has extensional
                                       files, which are interpreted
                                       directly.
The file above was created with...

$ truncate 922 . . . . . 000 f

So how does grep work so quickly?

If it is an intensional file, then it just reads that.
Hence, in this case grep quickly realizes that there is no x and returns almost immediately.

Aside....

This is useful for virus scanners to exploit. It wants to skip parts that cannot possibly be viruses. That way, it catches them faster.

## Problem Areas in OS Design

- Incommensurate scaling

not everything scales at the same rate

cost/unit ↓ — economies of scale (Adam Smith, pin-factory)
cost/unit ↑ — diseconomies of scale (star network)

Things break as you grow with diseconomies.
Economies can cause waste.

- Emergent properties

larger systems have properties that smaller ones didn't
  • (Tacoma Narrows bridge) ← resonance frequency
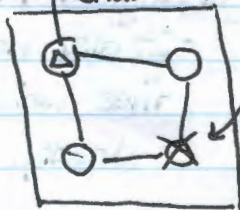  • UCLA campus network → Napster

- Propagation of effects

pathways for propagation are "more effective"

in a digital system, one component change can be very profound

international characters 天 ← \309\257



file ← C:\Windows\foo.txt should become 天.txt
system

Change Δ breaks part X.

Turns out, the coding was shift-JIS where first byte is all 1's, and the second could be anything.
However: the 2nd byte could be the ASCII '\', and the file system will not work even though it wasn't changed.

SOLUTION: use UTF-8.

- Trade Offs

 Waterbed Effect - if we push one part of the waterbed we must send the water somewhere else.
 [ex] if we used the SOLUTION, we would need to represent Japanese characters with ~~base~~ 3 characters.
 [ex] Too easy to break into UCLA registrar passwords too easy to guess
 SOLUTION: keycards
 Tradeoff: extra work for maintenance
 Benefit: Security.

## Lecture 2

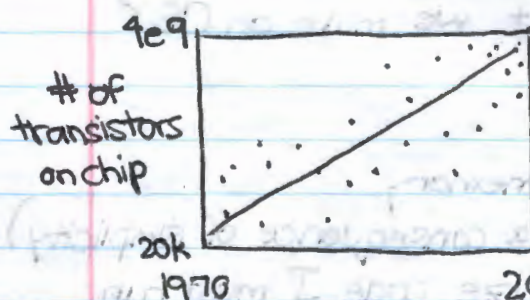1/6/2016

Today's topics:
- a bit more philosophy
- how not to make an OS

Complexity

that are economically viable

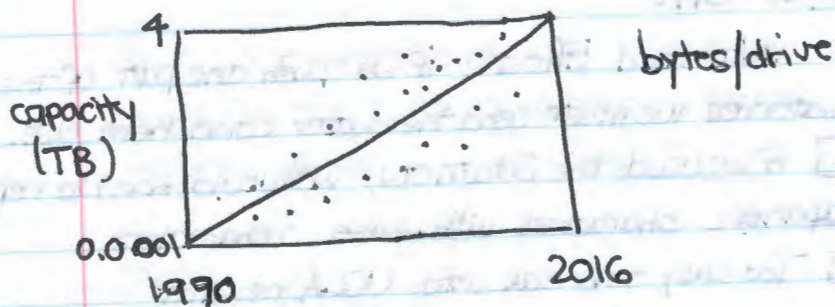Moore's Law - # of transistors doubles approx every year



Log scale! bits/chip

However, in recent years the exponential growth has begun to trail off.

# of transistors on chip

4e9

20k

1970          2016

Date of Introduction

Kryder's Law   secondary storage capacity

However, complexity is increasing at a higher rate

Why Exponential Growth?

→ UNIVAC I (slow, safe, hand designed)

use to design the UNIVAC II

UNIVAC II

$$\frac{d(technology)}{dt} = K * technology$$ , which comes out to be exponential.

## How Not to Make an OS

First off, Why <u>shouldn't</u> we make an OS?
- simplicity
- performance
  - speed   - memory
- Reliability (as a consequence of simplicity)
- security (minimize code I must run.

Consider that I am submitting a proposal for a business on Friday and I want to make my own system (I'm paranoid)

no OS,

Application — count of words in an ASCII text file,
`(bytes '\001' - '\177')

Interface
① Power switch

`1' - '127'          ~10 year old desktop

BIO's — where is program,     Core i3-4160 (3 MiB cache,
and where is text file?       4 GiB dual channel DDR3

39716

file is here → 1 TB hard drive, SATA, 7200rpm
in adjacent sectors  Intel 4400 graphics

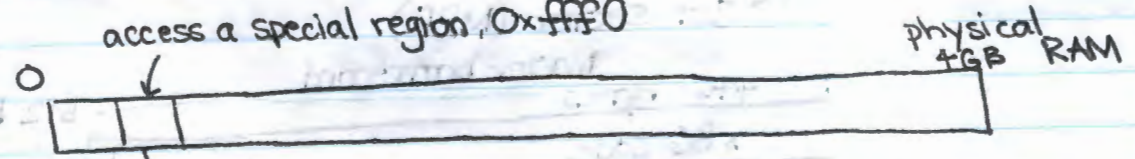Traditional disk drive, 512 bytes each sector

Cycling power clears CPU, cache, and RAM

Bootstrapping Problem — if we know ..., then......, but where
do we start? We need to have something to work with in
order to get started

built into hardware
instruction pointer — $0xffff0$     $2^{20} - 16$

access a special region, $0xfff0$                          physical
                                                           4GB RAM
0

ROM region ← hardwired by manufacturer to do
what we want it to do.
—or—
we program, if we got the manufacturer
PROM        to produce it for us. Most paranoid
(programmable approach.
read only
memory)
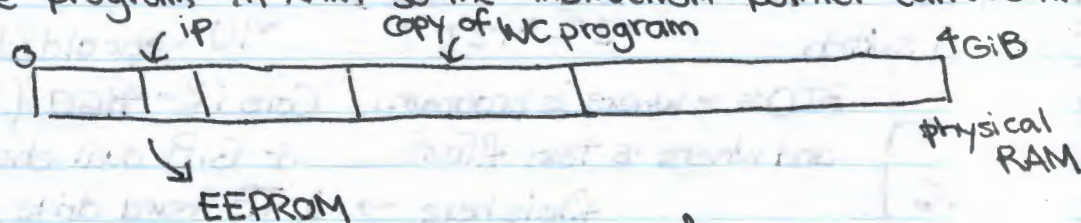—or—      called  "nonvolatile memory"

EEPROM
(electronically erasable    Downsides — might render computer
PROM)                        unbootable if bad erasure,
                             Erasure takes time and effort

We want the word count program to be on disk, and use the EEPROM to get to it somehow. The problem is, we need the program in RAM so the instruction pointer can use it.

ip          copy of WC program

0                                                           4 GiB

EEPROM

↳ location of wc *       + loading program
program, size of wc

Problem: * still
program dependent.     and so the solution is to have the
We don't want the     size and location to be the same for
EEPROM to have to     every program in EEPROM.
change.

Convention: (It's bad ☹)

                                                    disk

1st sector (MBR)
        master bootrecord
0      446 bytes                            512 B
        x86 code

copy to 0x7c00 in RAM,          convention:
then jump in.                   last two bytes
• hardware sanity checks        eventually have signature
  — some CPU checks  — checks RAM
• checks for devices in order      0x 55  0x AA
  finds 1st device with an MBR      little endian
                                     0x AA55
CPU
|              bus
    ①
devices

Firmware (EEPROM) → MBR (first sector on disk) →
    do whatever to code so long as it is in first 446 bytes

0    4  16                              64      2  512

Each entry will tell you a bit
about the rest of the disk
with sector counts

array of 4
16 byte items,
called the partition table, which tells
how the disk is divided.

MBR

6 ↑ origin + size

disk

32 bit    sector counts for offset partitions    size of partitions
other info: type byte, partition types, is Bootable?

<u>Chain Loading</u> - the firmware boots by reading the MBR, which
reads the first bootable ~~ports~~ partition (called a volume
boot record) and so on.

programs
→ more programs
→ even more programs

firmware → MBR → VBR → kernel
                ↑        ↑
         OS-agnostic  OS-specific
                        (many sectors)

We need a subroutine to read data from the disk.
needed by: programs, MBR, and the firmware

V₁   <u>cleaner in v2, pg 11.</u>              sector #    memory address

void read_ide_sector (int s, int *a)
                                  char*

Snippet 0  ✱  {
    while ((inb (0x 1f.7)       goes to disk controller to get data from
        & 0xc0) != 0x40)        disk controller, is slow
    continue;                   retrieve data from the bus address
    outb (0x 1f2, 1); ~~outb(0x1f~~
    outb (0x 1f3, s & 0xff);   || low order 8 bits
    outb (0x 1f4, s >> 8);   →  outb (0x 1f 6, s >> 24);   0x20 is
                                outb (0x 1f7, 0x20)        read

CPU

bus

disk controller

Computer

The controller communicates through the buses via their registers

$0 \times 1f7 \rightarrow$ status register 1B
contains info about state of disk's controller

| 0 | 1 | |

controller is available!

~~1f 7    status, command~~
~~1f 2    sector count~~
1f0    read data
1f2    sector count
1f3    — low order byte
1f4
Sector # —  1f5
1f6    — high order byte
1f7    status, cmd

```
* While (( inb(0 × 1f7) & 0x0 512 != 0x40
              continue;          ← sizeof(int)
      insl (0x1f0, a, 128);
   }
```

✓ We can have a copy in firmware at some location
            —or—
✗ it can be in the master boot record
            —or—
✗ it can be in the word count program

Isn't this wasteful? Let's just use one copy

BIOS = basic input/output system
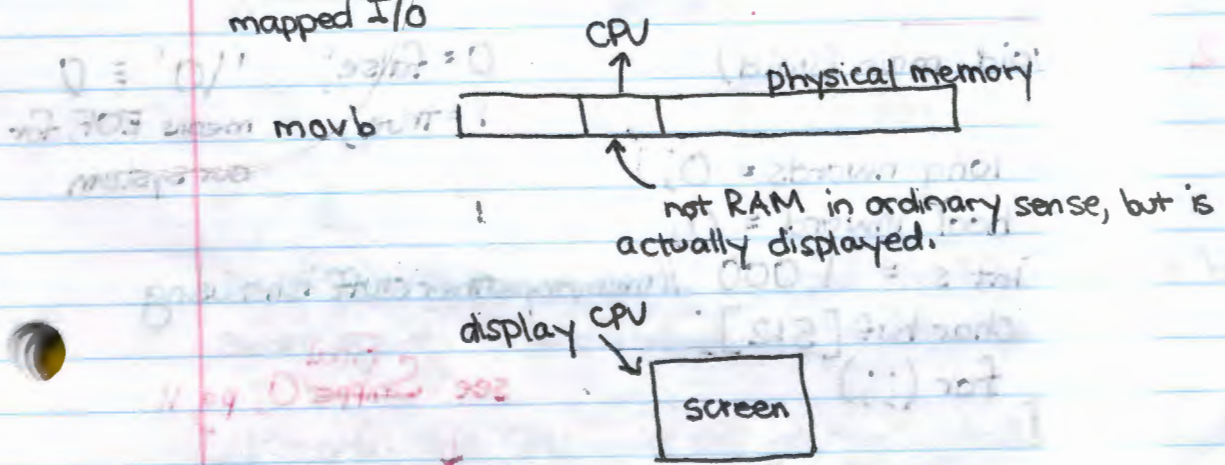downside: this must be built into the system

Final O

$V_2$

sector #    memory
address

```
void read_ide_sector (int s, char* a)
{                        4byte   4byte
    while ((inb( 0x1f7)& 0xc0) != 0x40)
        continue;
    // While not free, continue:
    outb ( 0x1f2, 1);    // Tells that sector count is just 1
    outb (0x 1f3, s);    // Write low 8 bits to s.
                                        from
    outb (0x 1f4, s >> 8); // Write next 8 bits to proper spot in s.
                                        from
    outb (0x 1f5, s >> 16); // Write next 8 bits to proper spot in s.
                                        from
    outb (0x 1f6, s >> 24); // Write next 8 bits to proper spot in s.
    outb (0x 1f7, 0x 20); // change status to reading the sector
    while ((inb(0x 1f7) & 0xc0) != 0x40)
        continue;
    insl ( 0x 1f0, a, 128); // Read 128 longs worth of memory into RAM
                    ↑            // starting from a, the mem address
            512 / sizeof (int)
}
```

(status & 11000000)
!= 01000000

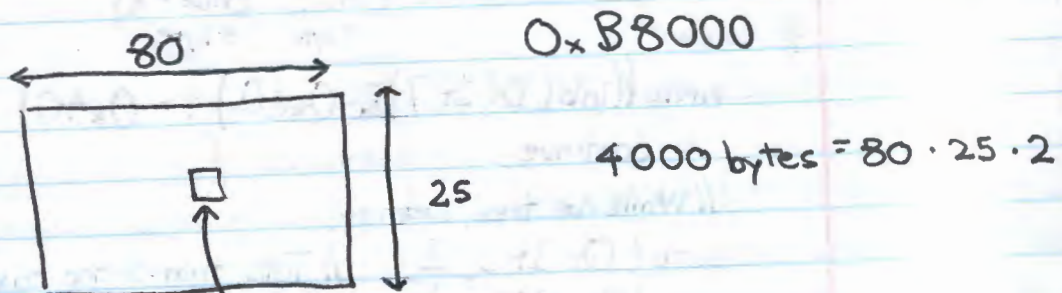01 . . . . . . . .
~~
sector is free

Lecture 3                                  1/11/16

To output to screen, not programmed I/O (PIO) but memory
mapped I/O

CPU
↑
physical memory
movb π  [ |      |      |      ]
         ↑
         not RAM in ordinary sense, but is
         actually displayed.

display CPU
    ↓
 [ screen ]

Standard:

Model screen as 80×25 grid of characters

$0 \times B8000$



4000 bytes = 80 · 25 · 2

25

16 bits to represent
low order 8 bits is ASCII
high order bits = graphical appearance

64 bit integer

$V_1$ : cleaner in v2 pg

**Snippet 1**

```
void output_to_screen (long n)          figure out where
{                                       onscreen to start
short   long *p = (long*) 0xB8000 + 80·25 - 80;
                                                2      2

    *p = x;    // Won't work!. Output will correspond
otherwise blank               to early bytes on screen.
screen if 0.
        do while (n != 0)        graphic              reversed
        {   *--p       0x700  setting                 order!
            *p++ = (7 << 8) | (n % 10 + '0');
            n /= 10;
        } while (n != 0);
}
```

$V_1$

**Snippet 2**

```
void main (void)                0 = false;    '\0' ≡ 0
{                               1 = true;   means EOF, for
                                            our system
    long nwords = 0;
    bool inword = 0;
    int s = 1,000   // memory other stuff isn't using
    char buf [512];                    Final
    for (;;)                    see Snippet 0, pg 11
    {
```
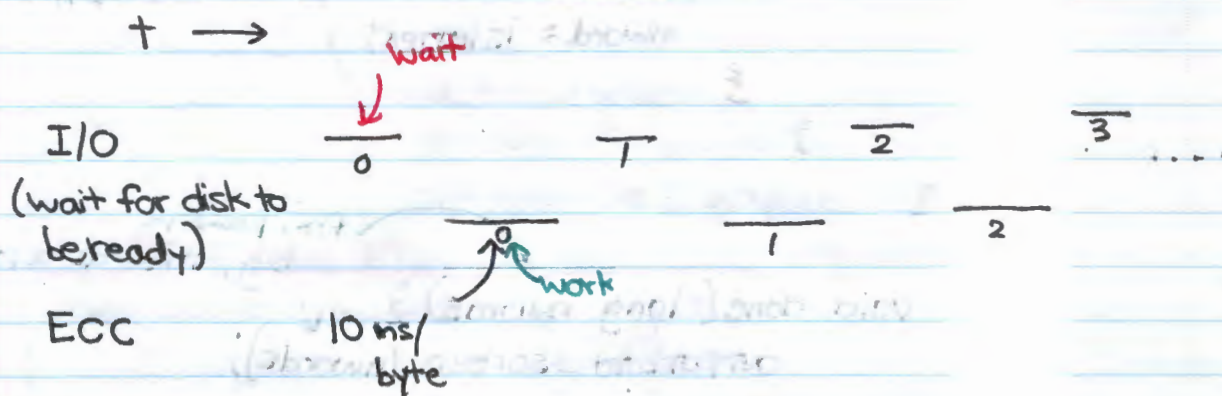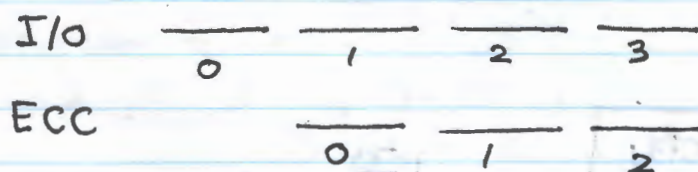
This adds complexity!. Furthermore, since we're looking at the data in RAM anyway, it is sent to the CPU anyway, so DMA isn't as useful.
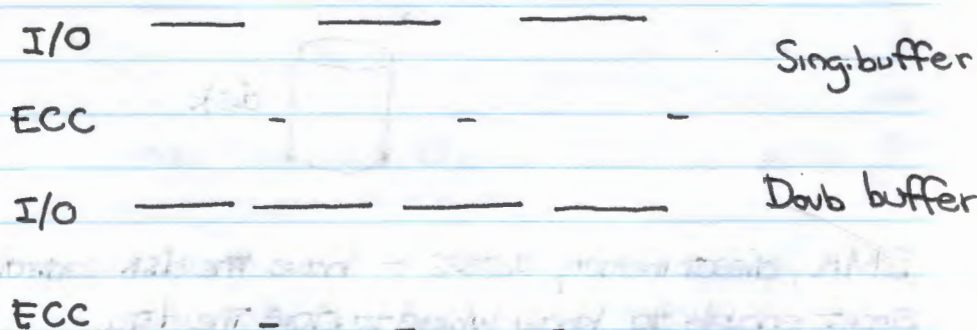
However, consider a crypto app; very CPU intensive

$t \longrightarrow$

I/O
(wait for disk to
be ready)

ECC     10 ns/
byte

How do we improve this? <u>Double buffering</u>, where we load one buffer while working on the other one.

I/O

ECC

However, if I/O and computation is too high in discrepancy using double buffering is not as worthwhile.

I/O                                    Sing. buffer

ECC

I/O                                    Doub buffer

ECC

Triple buffering is good when we load; process 1, process 2, etc.!

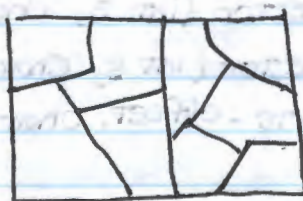Multitasking (with multiple processors)

I/O ≡≡≡   ≡≡≡

ECC

Load multiple from I/O then process one at a time. While loading next set, if processing is faster by a significant margin.

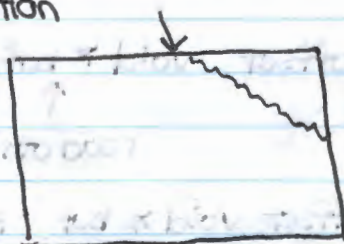How to scale up these programs.
- Fancier performance tricks
  without rewriting applications
- Multitasking
  without rewriting applications

Modularity

cost of maintaining a module of N lines of code is $O(N^2)$

Abstraction

Finding natural divisions in the programs that makes the program is easier to manage.

How do you measure the quality of modularity + abstraction?
Simplicity (ease of use)
           (ease of learning)
Robustness (tolerance of errors, large inputs) harsh conds
Performance   modularity costs this, minor costs unavoidable,
           avoid major costs as often as possible.
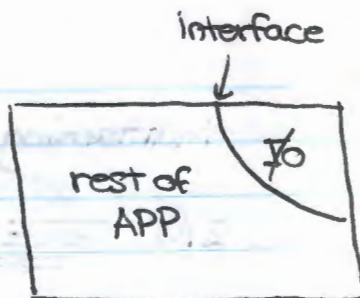Flexibility/Lack of Assumptions/Neutrality

Consider....

char * readline (FILE *f);
    BAD DESIGN for O.S.
    Why?
    ① Performance ✗ (unbounded work)
                    (unbatched work) ← overhead for small lines
    ② Robustness ✗ (apps crash for big lines)
    ③ Neutrality ✗ (forces particular line ending convention on the
                    app)
    ④ Simplicity ✓ 😄

Let's now examine read-ide-sector...


Interface

rest of
APP

I/O

void read-ide-sector (int s, char * buf)

improve
robust
ness by → int read-ide-sector (int s, char * buf)
reporting
potential      int read (int byte-offset, char * buf, int bufsize)
error                ↑
msgs.
            more general, no need to worry about exact sector size

        int read (int byte-offset, void * buf, int bufsize)
                                        ↑
                            read other things besides chars

            off-t
    int read (int byte-offset, void * buf, size-t bufsize)
            || long long on x86-64    ↑
                            increase max bufsize, make
                            this portable
                            || size-t is unsigned long on x86-64

    ssize-t read (off-t byte-offset, void * buf, size-t bufsize)
        ↑

    # bytes read        Other improvements....
    or -1 if error      have it return error code? Not the
                            convention, sadly.

which file? (opaque file handle)    where did the byte-offset go?

- ssize_t read (int fd, void * buf, size_t bufsize);

BIG IDEA OF UNIX
everything outside the program is a file
(disk, flash drive, network, mouse, keyboard, display)

random access          |   stream device
device                 |        ↑
                       |   doesn't need/use the offset; hence its
                       |   removal

but we lost random access! How do we get it back?

- size_st lseek (int fd, off_t ~~value~~ where, int flag);

Corollary, the OS records current file positions
starts at 0, after reading n we add n.
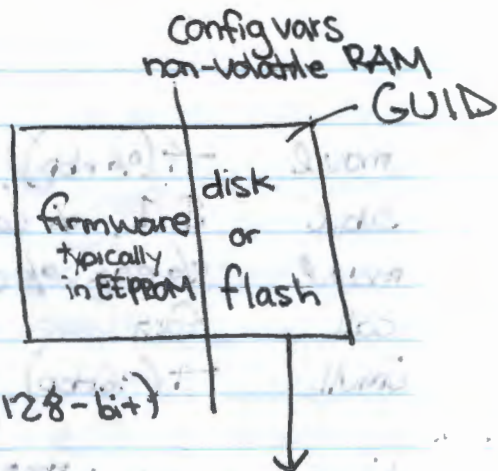
- pread ...
  add complexity to get performance

Mechanisms for Modularity
- Function calls
  - + simple and well understood
  + reasonably fast
  - things can go wrong
- Client-server (works, but.... slow)
- Virtualization

Lecture 4:
OS Organization

Booting UEFI

Config vars
non-volatile RAM

1/13/16

GUID



firmware
typically
in EEPROM

disk
or
flash

GUID: globally
unique ID for (128-bit)
disk partitions

Standardized — GUID partition table[++]
(GPT)

Standard layout in each partition:
  FAT 12      kernel with known name
  FAT 16
  FAT 32

Firmware is in charge until booting. This is <u>enforced</u> organization.

How to enforce modularity after booting   char buf[2000];
   function calls (terrible way?)   read(3, buf, 1000);

```
int fact (int n)               unoptimized:
{
    if (!n)                    fact:
        return 1;              push   pushq   %rbp
    return n * fact (n-1);   —( movq    %rsp, %rbp
}                                subq    $16, %rsp
                                 movl    %edi, -4(%rbp)
                                 cmpl    $0, -4(%rbp)
Recall, %edi is parameter        jne     .L2
        %rax is return value     movl    $1, %eax
                                 jmp     .L3
```

.L2

```
movl    -4(%rbp), %eax
subl    $1, %eax
movl    %eax, %edi
call    fact
imull   -4(%rbp), %eax
```
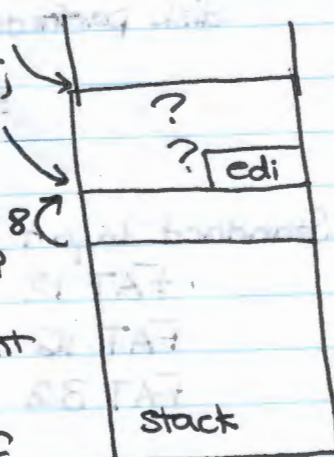
.L3

```
leave   rsp=rbp; rbp=*rsp++;
ret     rip = *rsp++;
```

pop = {

From start to finish:
- Decrement by 8
- Set %rbp to %rsp's current location
- Move up 16 bytes, 2 sets of junk inside.
- low order 4 bytes of %rdi is set to where %rbp is
- Perform comparison, if %edi is 0, go to L3, store 1 in %eax (low 4 bytes of %rax), then return

- Otherwise we jump to L2, first two commands get n-1 into %eax
- movl %eax, %edi moves, %edi
- Then we call factorial
- We then multiply into our %eax register which will contain the answer.
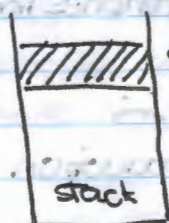
Things that can go wrong
fact ( INT_MIN) ⟵ overflow, attempt fact ( INT_MAX)
fact (50) ⟵ will only give us the bottom 32 bits (4 bytes)
of the correct answer

However, each time we call it we will put up 12 bytes,
eventually obliterating the rest of the data structures
(or dump core, if we're lucky)

We can catch this but they have some downsides...
· Adding checks will greatly slow it down
· Guard page



← if it reaches here we will trap

Suppose in the earlier machine code we call fact2 instead of
fact. What can fact2 do to mess up fact?
· Modify stored values in fact (activation record)
· make it return to the wrong spot
· modify the return value
· Cause an infinite loop.
· Jump elsewhere in the program
Can fact mess up fact2?
· Put a random value into %edi
· movl $0, %esp ⟵ will have fact2 enter
jump fact 2             forbidden zone as soon as
                        it tries to push to the stack

We have soft modularity.
- it doesn't scale to large applications

We want _hard modularity_, where a failure in a module won't tank the rest of them.
- virtualization     Unidirectional hard modularity
- client-server     Hard modularity in both directions

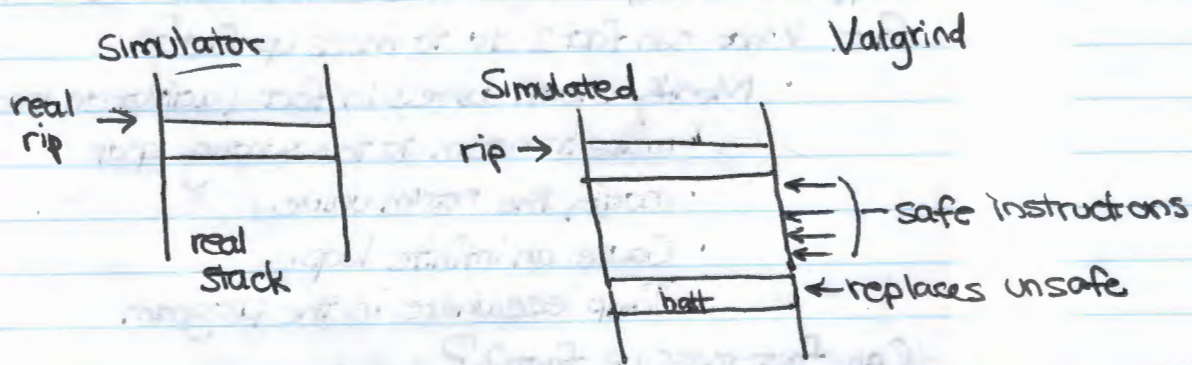_Simplest way to get virtualization_
   write a simulator for the machine the app runs on
      carefully checks instructions
        - halt → return
        - outofrangeaddr → return
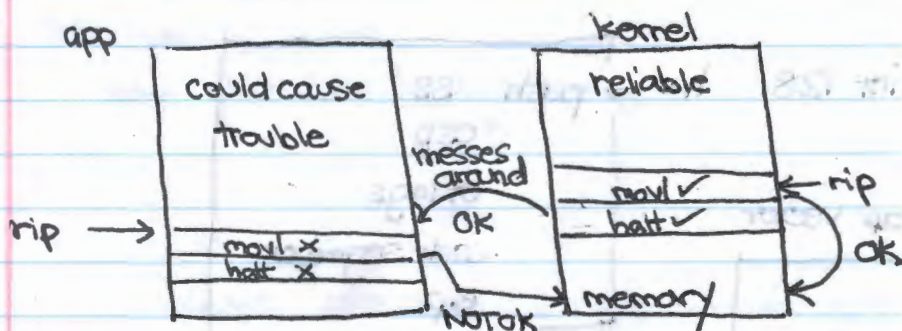      count instructions
        - When it reaches limit → return

Unfortunately, this method is slow. :(
In fact, it's usually too slow for production

We need a virtualizable machine, a virtualizable processor



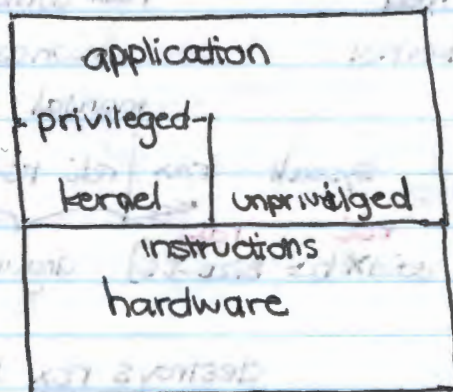There's one big problem left... it's still too slow!!!

app / kernel diagram:
- app: could cause trouble
- rip → movl ✗ / halt ✗
- messes around / OK / Notok
- kernel: reliable / movl ✓ / halt ✓ ← rip / OK
- memory

Protected Transfer of Control   ← should make these rare
for this model to work

~~dangerous~~ privileged instructions in the application
do not execute.

rip becomes in kernel

We have a privileged bit, which is set to 0 in the app.
When it encounters a privileged instruction, it passes to
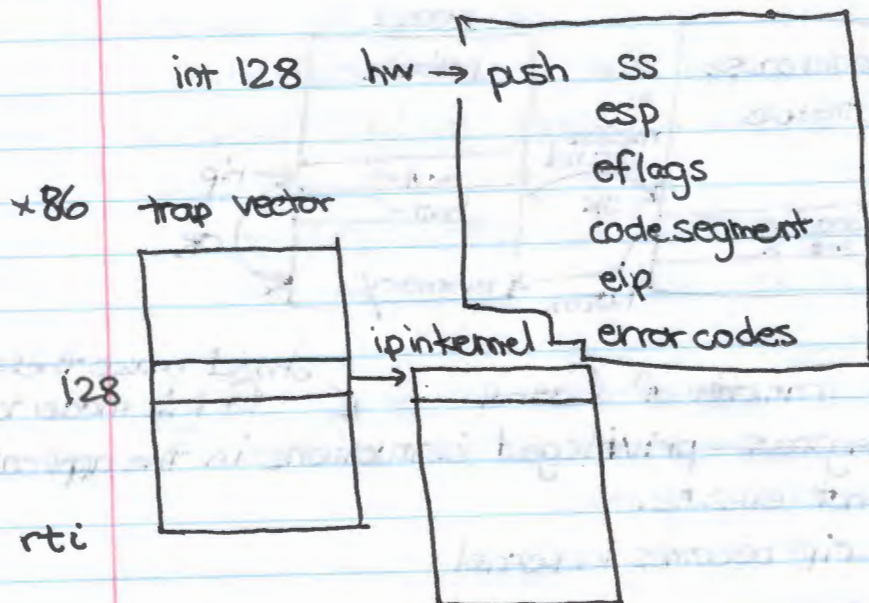the kernel, sets the privileged bit to 1, and does what
it needs to do.

Layered architecture



- application
- privileged →
- kernel | unprivileged
- instructions
- hardware

Classic way to enter the kernel:
   execute a privileged instruction with a standard convention
   this is → int 128

int 128    hw → push   SS
esp
eflags
code segment
eip
ip in kernel ⤷ error codes

x86   trap vector

128

rti

rti : int   as   ret : call
expensive   or   cheap

Nowadays, for the x86, x86-64, we have a special
sys call  _machine instruction_ :
- When executed
  - enters kernel
  - 

rax contains syscall #'s,
which can be found in a
manual

syscall   rax | rdi rsi rdx r10 r8, r9

rdi   rsi   rdx
ssize_t read(int fd, void* buf, size_t s)   arguments
{
   asm;
   deal with errno;          destroys rcx, r11
}                            result into rax
                            -4095  -1 means failure -errno

So if rax is wrong then we
send result to errno and return -1.

Mechanism works

Now, what does the user see?

↑

app
dev

Model: processes : programs running in a virtual machine atop an operating system

Creation   pid_t fork (void);   Clones current process
returns 0 in child
returns child's pid in parent
returns -1 on failure

Destruction   void _Noreturn _exit (int);

↑           ↑

never       bypass normal
returns     cleanup

Calling process immediately stops running, but not gone!

Process object          So this isn't the true way to
destroy a process

| exit_status | pid |
| fd-table | |

            status    flags — ∅
                ↙         ↙      ↘NOHANG
pid_t
~~int~~ waitpid (pid_t , int* , int )