# History of x86

4004
↓

8-bit | 8080 | 6502 | 6800

8085

Z80
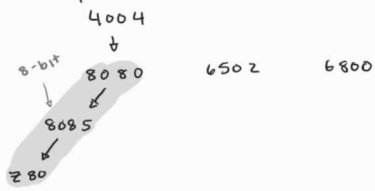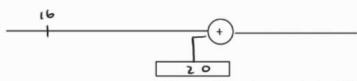
16-bit | 8086 | toy things - missing a lot
ex. no paging

32-bit | 80386

---

16 bit:

can address $2^{16}$ bytes ⇒ 64 KB
NOT A lot of space

you could buy 1MB chips but... there was
no way to address it —▷ 16 bits not enough
need segmentation



20 bits to point the address somewhere
↳ limit the amt of space for each
program but can run multiple
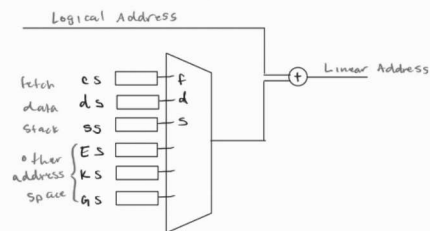programs
↳ how do we use 20 bits in a
16-bit system?
↳ assume lower 4 bits are 0
• each program gets its own 64 kb memory space
↳ useful if you want to have multiple programs at
a time but... no one really did that
↳ pretty much all compers were single user

What if we have 2 registers?
use a different address for fetch and everything
else
↳ have different section for program instructions
double the and data —▷ each 64 bytes(?)
address space

then... we can add more for
• stack
• other things...

Logical Address



fetch cs f
data ds d
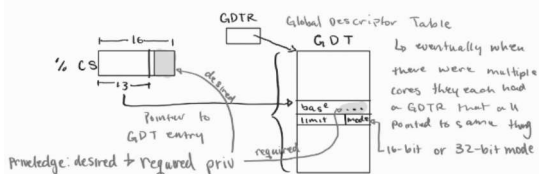stack ss s
other ES
address KS
space GS

• then there were too many and it got too crazy
✷ Note: you couldn't use C in this system

Then technology improved and we got 32 bit Addresses
At this point they could've just used paging
but... they didn't
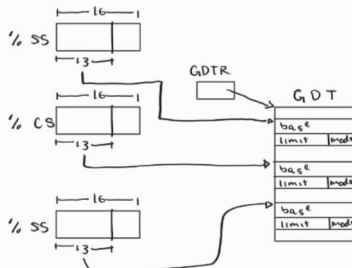
How do we widen the 16 bit value?
• want to keep 16 bits for backwards compatibility
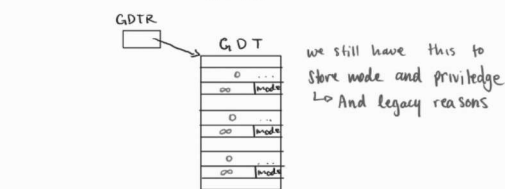• they... added a level of indirection



GDTR | Global Descriptor Table
% CS | GDT
↳ eventually when
there were multiple
cores they each had
a GDTR that all
pointed to same thing
Pointer to
GDT entry
Privledge: desired + required priv
base limit mode
required
16-bit or 32-bit mode

---

Eventually... created Local Descriptor Table
to have a GDT per process
• how do we know which DT to use?

% CS
↳ tells which table to use

This was **DUMB**. we don't have LDT anymore



% SS

% CS — GDTR — GDT
base limit mode
base limit mode
% SS
base limit mode

## NOW:

GDTR — GDT

we still have this to
store mode and privledge
↳ And legacy reasons



Logical Address

fetch cs f
data ds d
stack ss s

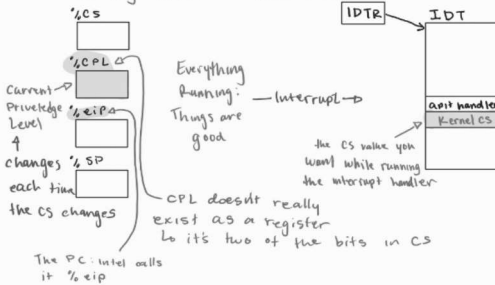+ → Linear Address → TLB/PD/PT → Physical Address

✷ now... if you run in 64-bit mode you
will get an exception if base and limit
are not 0 and ∞

---

## USER MODE

we have 2 values for CS:
• user cs } Points to entries in GDT
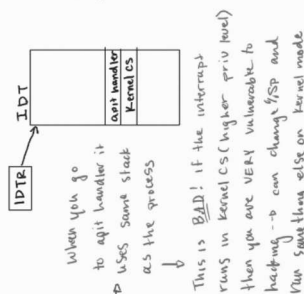  Both have 0 and ∞ for base and limit
• kernel cs } but mode is different

When a user program is running...
• CS reg MUST == user cs



% CS
% CPL
Current → Priveledge Level
4 changes %eip
each time
the CS changes %SP

IDTR — IDT

Everything
Running:
Things are —Interrupt—▷
good

cpt handler
Kernel CS

the CS value you
want while running
the interrupt handler

CPL doesn't really
exist as a register
↳ it's two of the bits in CS

The PC: intel calls
it %eip



IDT
cpt handler
Kernel CS

IDTR

—Interrupt—▷ When you go
to cpt handler it
uses same stack
as the process

This is **BAD**! if the interrupt
runs in kernel CS (higher priv level)
then you are VERY vulnerable to
hacking —▷ it can change %SP and
Run something else on kernel mode

Whenever you switch to higher privledge level
you MUST switch stacks
↳ there's a kernel stack!

How do we find kernel stack?
↳ logically you should have a reg that points
to kernel stack but... LEGACY

%CS
%CPL
%eip
%SP

PC
CS
flags

TSS ← All empty...
except

TSSP

Pushes these
so it can
return after

Kernel
Stack
PC
CS
%eip
SS

Memory:



0x00···        0x8····        0xF···

kernel

• If a user process tries to access something in kernel memory then the MMU will fail
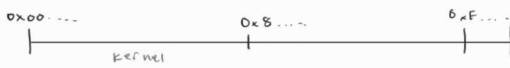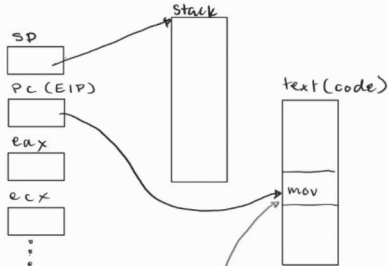  ↳ How does hardware know if it has access?
    ↓
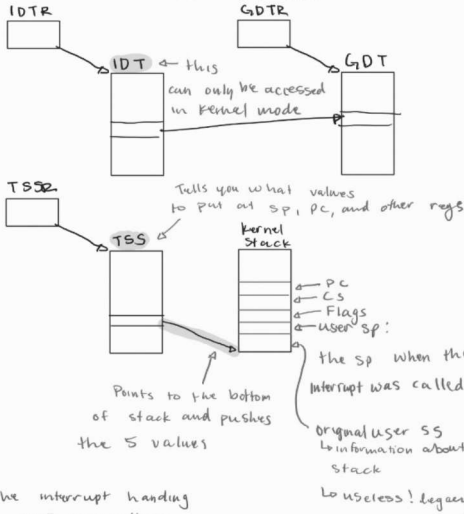  %CS contains CPL ──→ tells if user mode or kernel mode.

Lets start w/ user mode



SP
stack
PC (EIP)
text (code)
eax
mov
ecx
⋮

What happens if this creates a page fault
↓

• we should flush the pipeline ──→ this creates a big security risk ──→ easy to exploit
• need to run the interrupt handler



IDTR        GDTR
IDT ← this can only be accessed in kernel mode
GDT

TSS2        Tells you what values to put at SP, PC, and other regs
TSS
kernel Stack
← PC
← CS
← Flags
← user SP:
Points to the bottom of stack and pushes the 5 values

the SP when the interrupt was called

Original user SS
↳ information about stack
↳ useless! legacy

Note: the interrupt handling is NOT instructions.
     It is in hardware

✳ if we switch stacks we push 5 regs
✳ if we DON'T switch stacks we only push 3 regs
✳ SOME interrupts push an extra 6th reg with the error code
  ↳ It is the responsibility of the handler to take care of the extra reg
    ↳ iret doesn't know to take the extra reg off.
✳ Once the hardware finishes:
  - pushing the 5 or 6 regs
  - sets the SP
  - sets the PC
  - modify the flags as needed
• then you can context switch to the handler
  • change the CPL to 0 ──→ kernel mode
  • the handler pops things off the stack
    ↳ like the extra error reg
  • runs the handler
  • runs iret
• Returns to PC
  ↳ value to PC depends
    - for page fault PC = the instruction that failed
      ↳ allows the instruction to be run again
    - for other things the PC = the instruction after the failed instruction
      ↳ allows it to continue

○ If running in user mode:
  • There may be a lot of things I want to do that can only be done in kernel mode
    ↳ ex. allocating memory
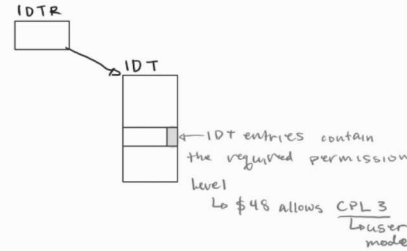      ↳ need access to TCBs and stuff which is only available in kernel mode

Solution?
use fake interrupt

        int $48 ←
  fake interrupt instruction

• tells hardware to act as if that interrupt is called



IDTR
IDT
← IDT entries contain the required permission Level
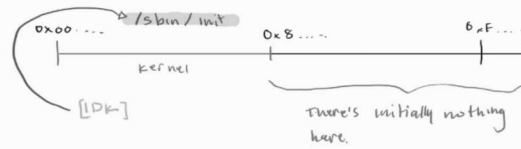  ↳ $48 allows CPL 3
      ‾‾‾‾‾
      ↳ user mode

• This still runs everything from before
  ↳ pushing to kernel stack & stuff

• you still need to tell it what to do
  ↳ use %eax and put a number
    ↳ system call: software interrupts

    ↳ different numbers mean different things
      ↳ the handler checks if you're allowed to perform the operation

Now. How do I enter User mode in the first place?



0x00···  /sbin/init   0x8····        0xF····
kernel
[IDK]       There's initially nothing here.

• kernel main shouldn't be running silly programs
  ↳ it should be loading/running user code
• kernel needs to mmap the user program
  ↳ need to set up the world:
    - set CPL to 3
    - set the SP
      ↳ switches the SP but you're still running kernel code but with a user stack
    - set PC

remember:
there's no instruction to set CPL
↳ the only two ways are exception handler and iret

• How do we set these 3 things?
  ↳ iret conveniently does these 3 things
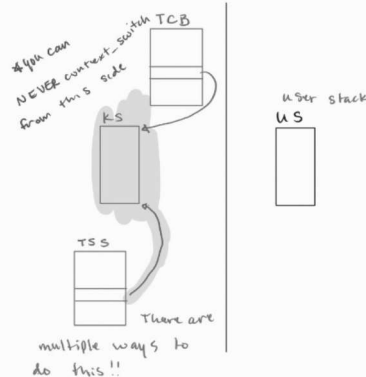      ⇓
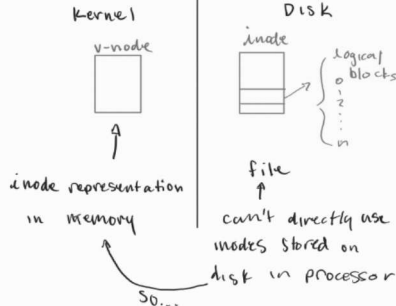    push 5 values on the stack
    then call iret

What happens to the Kernel Thread?
• It's still there! Kernel Thread is the thing that called iret!
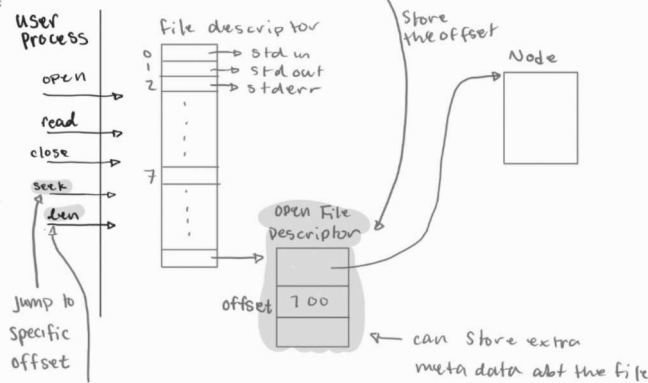• The kernel Thread has "morphed" into a user process
• When an exception happens... you're still on the same thread!



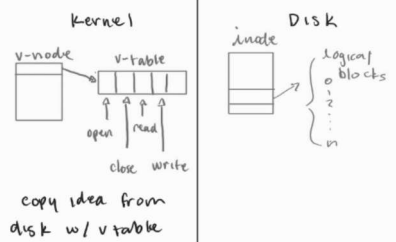✳ you can NEVER context-switch from this side

TCB
KS
user stack
US
TSS

There are multiple ways to do this!!

user process    Kernel    Disk

v-node    inode

logical blocks 0 1 2 ... n

inode representation
in memory

file

can't directly use inodes stored on disk in processor

So...

what if ... I wanted my file system to be not local for my machine
└ ex. labmachines

---

user process    Kernel    Disk

v-node  v-table    inode

logical blocks 0 1 2 ... n

open  read
close  write

copy idea from disk w/ v table

---

user process    Kernel    Disk

open

make call to kernel which calls to disk to get file ... now what?

Strongptr  Node    inode

logical blocks 0 1 2 ... n

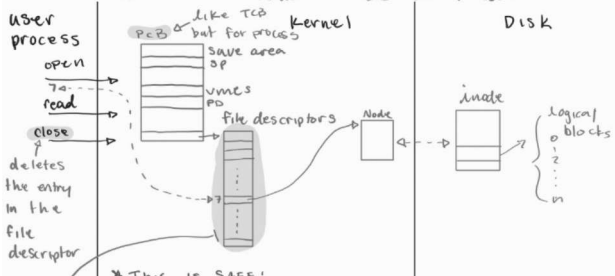you tend to want to reference count your access to node
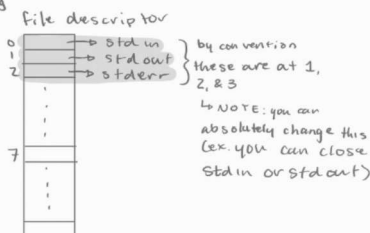└ coordinate access to Node

what about access in user process?
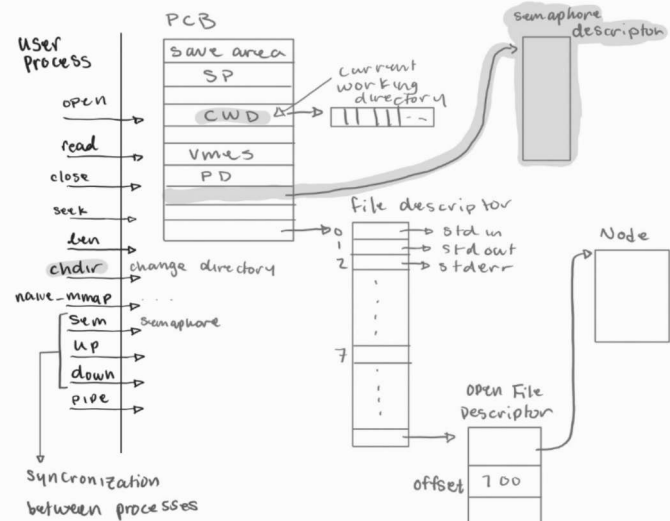└ can't directly read from disk
  └ need to use system call

open() returns the pointer to the memory
└ the user process shouldn't have access to the file cause it's stored in kernel heap
  └ what if we only let user process    └ Protected call read() using the pointe?
    └ how do you know if its a pointer you gave user? --> you don't. HUGE security risk
  └ *open() CANNOT return a direct pointer --> too unsafe
  └ use indirection --> open returns a pointer to the pointer to the file
    └ NEVER leak kernel addresses out to user

---

user process    Kernel    Disk

open
70--
read
close

PCB  like TCB but for process
save area
SP
vmes
PD
file descriptors    Node    inode

logical blocks 0 1 2 ... n

deletes the entry in the file descriptor

*This IS SAFE!
└ if you try and use a file desriptor value that WASNT given to you, either:
1) the index in file descriptor is invalid and creates an error
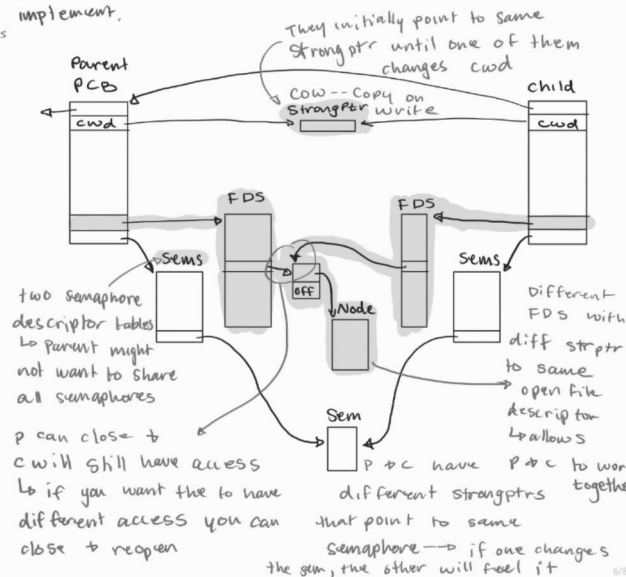2) it will give you some other file that you already had access to.

file descriptor
0  std in
1  std out
2  std err    } by convention these are at 1, 2, & 3
7            └ NOTE: you can absolutely change this (ex. you can close std in or std out)

---

Remember most of the time we read, we are going to read multiple chunks in order

user process

open
read
close
seek
len

file descriptor
0  std in
1  std out
2  std err

7

Store the offset

Node

Open File Descriptor

offset  700

└ can store extra meta data abt the file

Jump to specific offset

in our prog:
simplification of 'stat' command
└ length of file => size in bytes

---

user process

open
read
close
seek
len
chdir  change directory
name_mmap
sem  Semaphore
up
down
pipe

PCB
save area
SP
CWD
vmes
PD

current working directory

semaphore descriptor

file descriptor
0  std in
1  std out
2  std err

7

Open File Descriptor

offset  700

Node

Syncronization between processes

---

Fork! 

int id = fork();    // makes a clone of the process
                    └ they have the same history

└ returns int:
  0 => child process
  >0 => Parent -- returned value is id of child
    └ allows you to do things like wait for child
  <0 => Error: you are parent and no child was created

Parent knows children but children don't know parents
└ there is system call for child to get parents but WE will not implement.

They initially point to same Strong ptr until one of them changes cwd

COW -- Copy on Strongptr write

Parent PCB    child
cwd    cwd

FDS    FDS

Sems    off    Node    Sems

two semaphore descriptor tables
└ Parent might not want to share all semaphores

P can close + C will still have access
└ if you want the to have different access you can close + reopen

Sem

P & C have different strongptrs that point to same semaphore --> if one changes the sem, the other will feel it

Different FDS with diff strptr to same open file descriptor
└ allows P & C to work together

What about File system?

Parent PCB

Deep copy:
copy everything

child



copy

copy

very expensive operation
So...

Parent PCB

child



Read only pointer to
same file until...
  ↳ EITHER Parent
OR child writes
  ⇕
then make copy of file
  ↳ COW → copy on write

Wait:
• parent waiting for child to exit
  ↳ returns the exit code
• you can disown a child
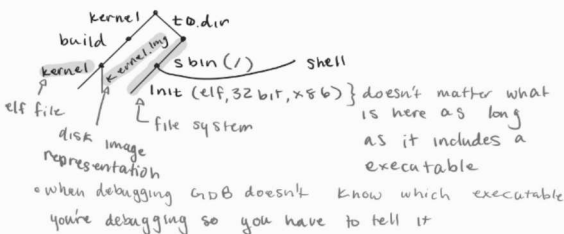  ↳ close()
  ↳ tells kernel to clean after the child.

# Welcome To User Mode

now:
• We REALLY want to avoid panicing
  ↳ only panic when we REALLY don't know what to do
• Tests no longer have 1 kernel per test
  ↳ it doesn't recompile for each test
    ↳ it just reboots
    ↳ unless you make clean

test cases structure:



kernel  t0.dir
build
kernel  kerneling  sbin(/)  shell
Init (elf, 32 bit, x86) } doesn't matter what
is here as long
as it includes a
executable

elf file
disk image
representation

file system

• when debugging GDB doesn't know which executable
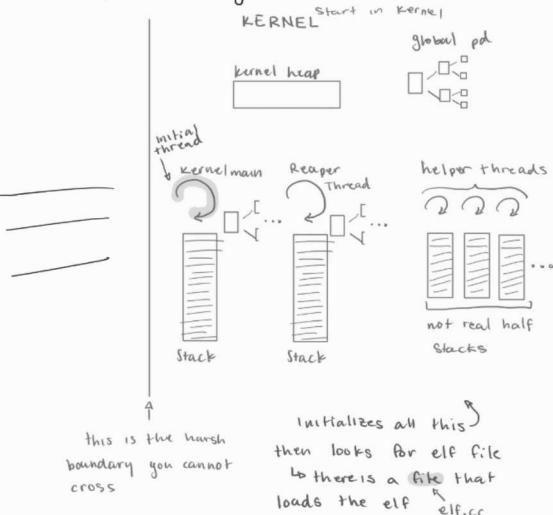you're debugging so you have to tell it

Notes:
• There's no printf in the test case → you have to move
some printf functionality to tc and do a write
system call
So what do we need in tc?
- dir directory
- user process has its own memory space
- syscall.cc / syscall.h has the system calls
  ↳ Missing 3 -- Described in README in t0

what happens when you run t0?

KERNEL  Start in Kernel
global pd

kernel heap



initial thread
kernelmain  Reaper Thread  helper threads

Stack  Stack  not real half stacks

this is the harsh boundary you cannot cross

Initializes all this
then looks for elf file
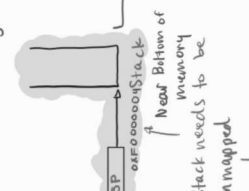  ↳ there is a file that
loads the elf  elf.cc

What is the elf file?
• has information on the data of the executable
  ↳ ex how big, where it is, which va's to use, etc
• The elf file can not handle position independent
• It's a good idea to demand page the program
  ↳ most programs will not run every instruction
• The elf loader is NOT lazy → might be a
good idea to change it to be lazy
• You will need to modify your mmap implementation
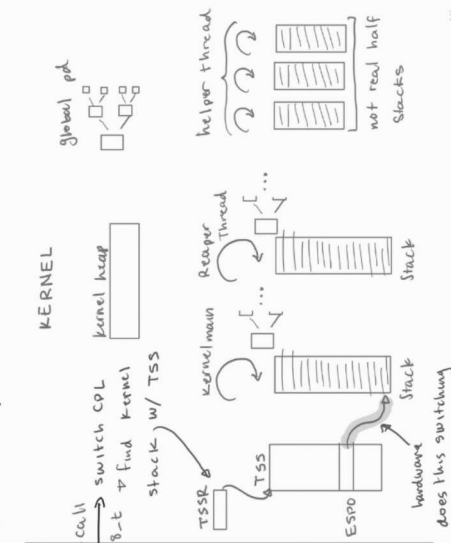to take an optional parameter for va—

If we allocate stack lazily
how will the page fault know
it's from stack?
  ↳ you can store the location of
stack in TCB
  ↳ in page fault it's close
enough to stack the we just
give the page

# Note: There will be some upper limit
in stack size → there will be
an Ed post

User mode:
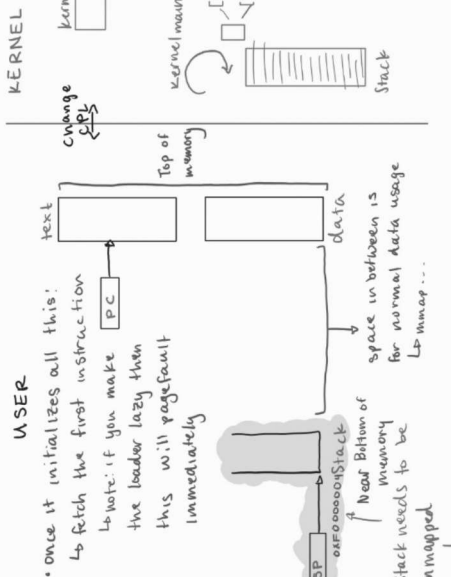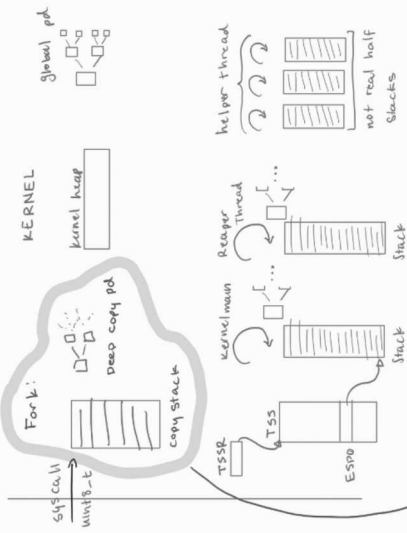• Ideally you want the stack to be far far away
from your program code
• Once it initializes all this:
  ↳ Fetch the first instruction
  ↳ Note: if you make
the loader lazy then
this will pagefault
immediately

KERNEL
  global pd
  kernel heap



helper thread  not real half stacks

Reaper Thread
kernelmain
Stack

Fork:
Deep copy pd



copy stack

TSSR  TSS  ESP0

syscall
uint8_t

Fork

minimum
hello love

• Parent will run on kernel stack
• child will run on its own stack
• the expensive part of the deep copy of pd

Exec:
execl (path, argv0, argv1, ..., 0)
- Stops current program and switches to
the program at path
- If successful → it will NEVER return
  ↳ you need to empty your address space
  ↳ retrieve the new elf file
  - create brand new stack
  - pass the arguments to new program
    ↳ need to load the args into kernel
stack before deleting old address space
    (↳ load back to new address space
  ↳ doesn't matter where in address space
as long as the stack points to
the right address.

KERNEL
  global pd
  kernel heap



helper thread  not real half stacks

Reaper Thread
kernelmain
Stack

Stack

When you do a system call:

syscall
uint8_t  switch CPL
→ find kernel
stack w/ TSS

TSSR  TSS  ESP0  hardware
does this switching

KERNEL
  kernel h...
  kernelmain
Stack

change CPL

Top of memory

text  PC

USER  data

space in between is
for normal data usage
  ↳ mmap...

0xF0000000 Stack

↳ Near Bottom of
memory

SP

Stack needs to be mmapped