

EECS-343 Operating Systems

Lecture 3:

Process Creation and Memory Layout

Steve Tarzia

Spring 2019

Northwestern

Announcements

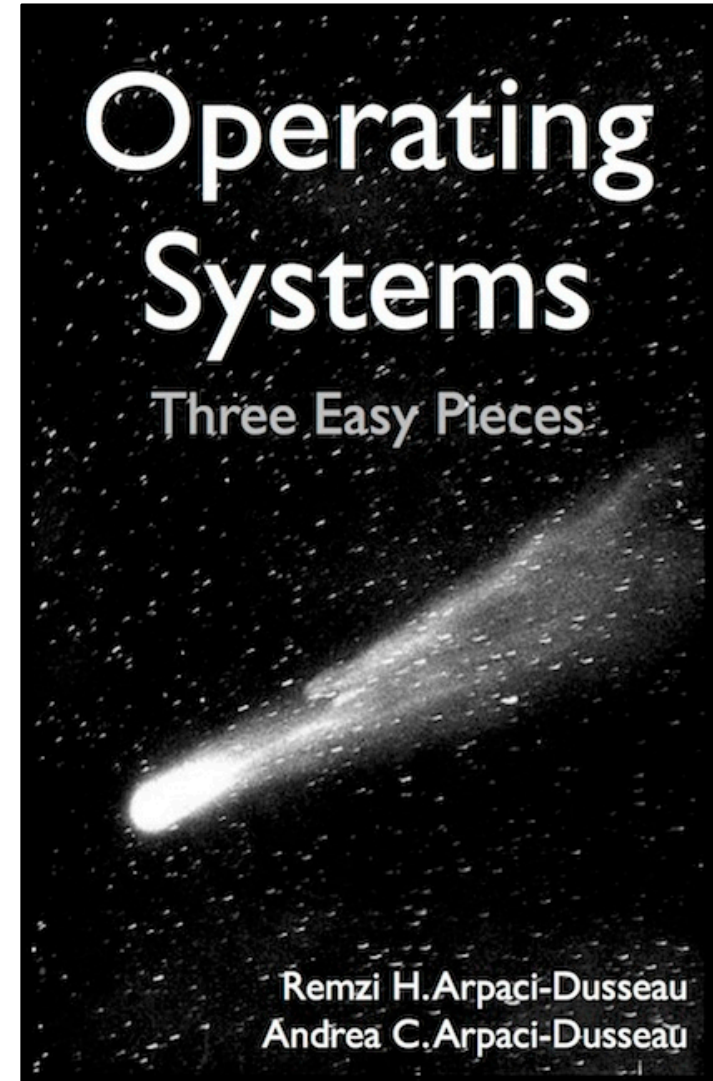
- Project 1 due on Monday!

Last Lecture

- **Process** is a program in execution
- **Limited direct execution** is a strategy whereby a process usually operates as if it has full use of the CPU & memory.
- CPUs have user and kernel **modes** to prevent user processes from running privileged instructions, thus *limiting* execution.
- **Interrupts** are events that cause the kernel to run
- **System Calls** (or traps) are software interrupts called by a user program to ask the OS to do something on its behalf.
- **Timer Interrupt** ensures that the kernel eventually runs.

Readings

- So far, we've covered Chapters 1-4 and 6 (Chapter 5 is today).
- Please read the Scheduling chapters next (Chapters 7-9)
- In the future, just try to follow along on your own.
- The syllabus says which chapters we're skipping.



Example Unix syscalls (process-related)

- `exit` – terminate the current process
- `fork` – duplicate the current process
- `wait` – wait for a process to terminate
- `exec` – run a program (in the current process)
- `time/stime` – get/set current time (in seconds)
- `brk` – change the process “break,” meaning max memory address
- `getpid` – get current process’s id
- `pause` – wait for a signal from another process
- `kill` – send a signal to another process (named after one signal type)
- `getuid/setuid` – get/set the effective user id of the current process

Example Unix syscalls (file-related)

- `read/write` – read/write data from a file descriptor
- `open` – open/create a file
- `close` – close a file descriptor
- `chdir` – change working directory
- `mknod` – create a filesystem folder
- `chmod` – change permissions of a file
- `chown` – change ownership of a file
- `seek` – change r/w offset in a file
- `utime` – change modification time of a file/folder
- `mount/umount` – mount or unmount a filesystem

“Hello world” with syscalls (in Linux)

C code:

```
int main() {
    write(1, "Hello, world\n", 13);
    exit(0);
}
```

- Notice that we are not using **printf**
 - printf is a libc function
 - libc’s implementation of printf will use **write**, which is a syscall.

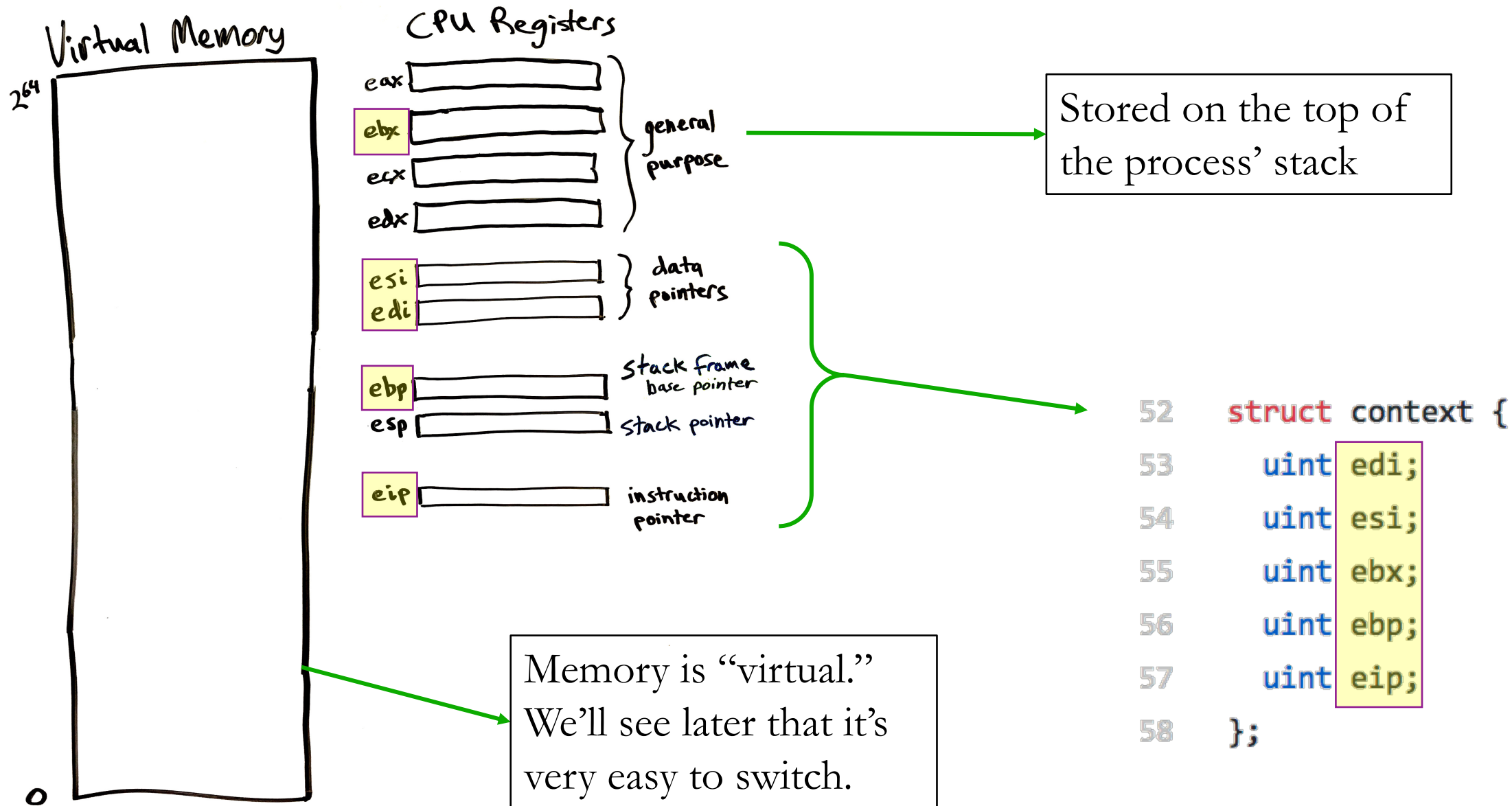
```
1  .section .data
2  string:
3      .ascii "hello, world\n"
4  string_end:
5      .equ len, string_end - string

6  .section .text
7  .globl main
8  main:
    First, call write(1, "hello, world\n", 13)
9      movl $4, %eax           System call number 4
10     movl $1, %ebx           stdout has descriptor 1
11     movl $string, %ecx      Hello world string
12     movl $len, %edx         String length
13     int $0x80               System call code

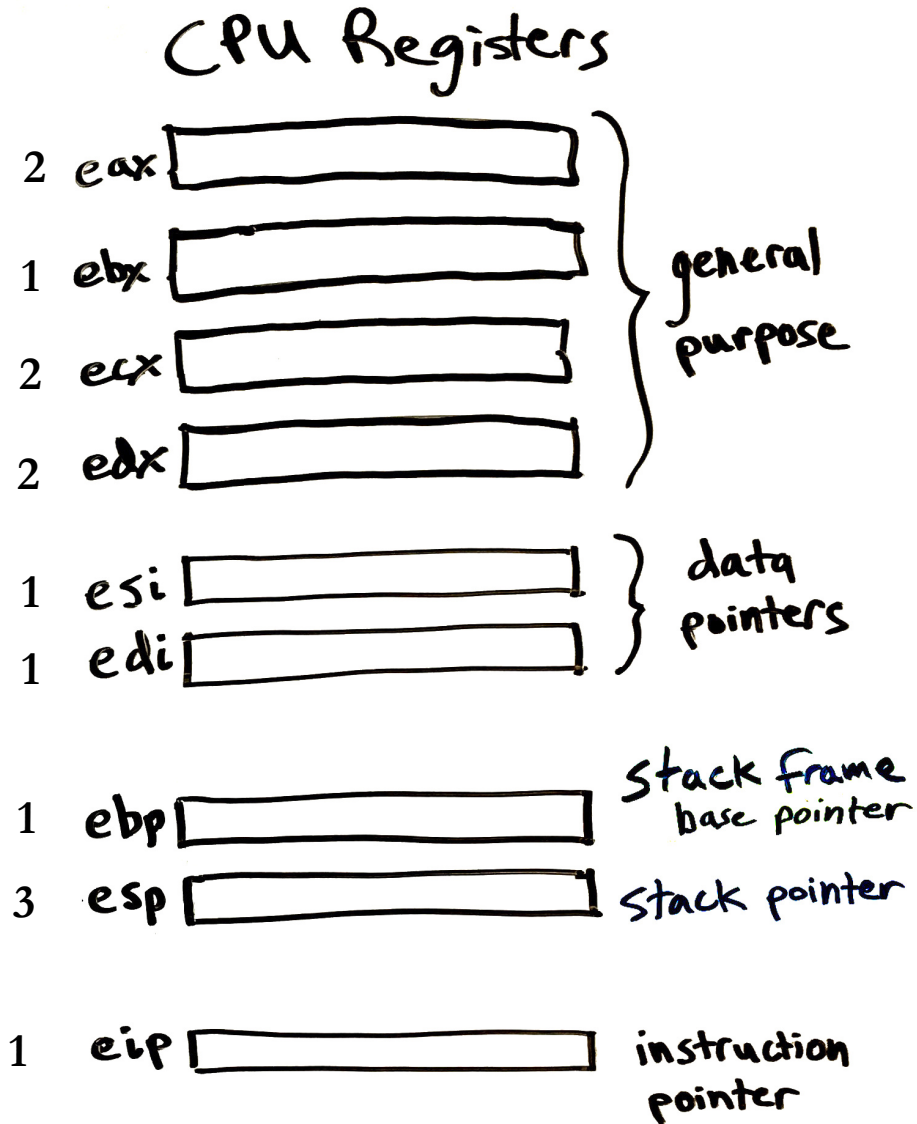
    Next, call exit(0)
14     movl $1, %eax           System call number 0
15     movl $0, %ebx          Argument is 0
16     int $0x80               System call code
```

(Bryant and O’Hallaron, Figure 8.11) →

Last time: Arrows on this slide were wrong



xv6 stores register values in *three* places



1. In **struct context** (`proc->context`):
`ebx`, `esi`, `edi`, `ebp`, `eip`
and `esp` is the address of the struct.
 2. In the user process' **stack**:
`eax`, `ecx`, `edx`
(by the x86 calling convention)
 3. In **struct trapframe** (`proc->tf`):
`esp`, *and also copies of*
`edi`, `esi`, `ebp`, `eax`, `ebx`, `ecx`, `edx`
 - These are automatically written by the CPU hardware when an interrupt occurs.
- Why store duplicates? ...*idk*

proc.h

```
17 // Saved registers for kernel context switches.
18 // Don't need to save all the segment registers (%cs, etc),
19 // because they are constant across kernel contexts.
20 // Don't need to save %eax, %ecx, %edx, because the
21 // x86 convention is that the caller has saved them.
22 // Contexts are stored at the bottom of the stack they
23 // describe; the stack pointer is the address of the context.
24 // The layout of the context matches the layout of the stack in swtch.S
25 // at the "Switch stacks" comment. Switch doesn't save eip explicitly,
26 // but it is on the stack and allocproc() manipulates it.
27 struct context {
28     uint edi;
29     uint esi;
30     uint ebx;
31     uint ebp;
32     uint eip;
33 };
```

Pushed on "top"

When kernel takes over during the interrupt handler, it copies register values from the trap frame to a new struct context that's pushed on the user process' stack.

...and in x86.h:

```
148 // Layout of the trap frame built on the stack by the
149 // hardware and by trapasm.S, and passed to trap().
150 struct trapframe {
151     // registers as pushed by pusha
152     uint edi;
153     uint esi;
154     uint ebp;
155     uint oesp;    // useless & ignored
156     uint ebx;
157     uint edx;
158     uint ecx;
159     uint eax;
160
161     // rest of trap frame
162     ushort gs;
163     ushort padding1;
164     ushort fs;
165     ushort padding2;
166     ushort es;
167     ushort padding3;
168     ushort ds;
169     ushort padding4;
170     uint trapno;
```

```
172     // below here defined by x86 hardware
173     uint err;
174     uint eip;
175     ushort cs;
176     ushort padding5;
177     uint eflags;
178
179     // below here only when crossing rings, such as from user to kernel
180     uint esp;
181     ushort ss;
182     ushort padding6;
183 };
```

The OS Coder's Curse



Photo from <http://www.righto.com/2013/09/intel-x86-documentation-has-more-pages.html>

- Not only do we have to use C...
- We also have to understand the Intel x86 processor architecture
- x86 is messy because it carries
 - *40 years* of incremental updates and backward compatibility
 - but it's the architecture most relevant to SW Eng. practice
- We'll gloss over some of the low-level details
 - Read the xv6 book & code when you really need to know.

Interrupt handling involves both hardware and software

In response to interrupt, the CPU *hardware*:

- Saves main registers to trap frame on the *kernel stack* (each process has two stacks)
- Switches to kernel mode
- Jumps to interrupt handler code

Then kernel *software* takes over to handle the interrupt and when finished can switch to a different process if desired.

OS @ boot (kernel mode)	Hardware	
initialize trap table	remember address of... syscall handler	
OS @ run (kernel mode)	Hardware	Program (user mode)
Create entry for process list Allocate memory for program Load program into memory Setup user stack with argv Fill kernel stack with reg/PC return-from-trap		
	restore regs from kernel stack move to user mode jump to main	Run main() ... Call system call trap into OS
	save regs to kernel stack move to kernel mode jump to trap handler	
Handle trap Do work of syscall return-from-trap		
	restore regs from kernel stack move to user mode jump to PC after trap	... return from main trap (via exit ())
Free memory of process Remove from process list		

Table 6.2: Limited Direction Execution Protocol

Instruction set architectures vary

- Low-level OS code for Intel x86 looks very different than that for ARM, PowerPC, SPARC, etc.
- Linux supports all of the above architectures and it requires different assembly code to handle context switches and interrupts on each.
- So, let's try not to get hung up on the machine-dependent details.



An OS can support multiple CPU architectures

- Linux supports x86 plus 30 other architectures, and growing!
 - See <https://github.com/torvalds/linux/tree/master/arch>
- How? Different low-level code is used for different builds.
 - Includes some C and Assembly code
 - This is just a small fraction of the overall Linux codebase
 - But it would probably be close to half of xv6, since it's such a simple OS.
- “Ports” of the Linux OS tend to be managed by different groups
 - Eg., much of the ARM source code bears the following comment:

```
Copyright (C) 2012 ARM Ltd.  
Authors: Will Deacon <will.deacon@arm.com>  
Catalin Marinas <catalin.marinas@arm.com>
```

Writing an OS for multiple CPU architectures

What's different?

- All the assembly code + some C
- Boot code
- *Mechanisms* for
 - Interrupt handling
 - Context switching
 - Memory management
- Device drivers (to control peripheral hardware)
- Etc.

What's the same? ... most C code:

- Filesystems
- Process scheduler
- Inter-process communication
- Networking
- Security / user management
- *Policies* for
 - Context switching
 - Memory management

Linux's entry.S in both x86 and arm for context switch

```
224  /*
225  * %eax: prev task
226  * %edx: next task
227  */
228  ENTRY(__switch_to_asm)
229      /*
230      * Save callee-saved registers
231      * This must match the order in struct inactive_task_frame
232      */
233      pushl   %ebp
234      pushl   %ebx
235      pushl   %edi
236      pushl   %esi
237
238      /* switch stack */
239      movl    %esp, TASK_threadsp(%eax)
240      movl    TASK_threadsp(%edx), %esp
241
242
243
244
245
246
247      /* restore callee-saved registers */
248      popl    %esi
249      popl    %edi
250      popl    %ebx
251      popl    %ebp
252
253      jmp     __switch_to
254  END(__switch_to_asm)
```

```
954  /*
955  * Register switch for AArch64. The callee-saved registers need to be saved
956  * and restored. On entry:
957  *   x0 = previous task_struct (must be preserved across the switch)
958  *   x1 = next task_struct
959  * Previous and next are guaranteed not to be the same.
960  *
961  */
962  ENTRY(cpu_switch_to)
963      mov     x10, #THREAD_CPU_CONTEXT
964      add     x8, x0, x10
965      mov     x9, sp
966      stp    x19, x20, [x8], #16      // store callee-saved registers
967      stp    x21, x22, [x8], #16
968      stp    x23, x24, [x8], #16
969      stp    x25, x26, [x8], #16
970      stp    x27, x28, [x8], #16
971      stp    x29, x9, [x8], #16
972      str    lr, [x8]
973      add     x8, x1, x10
974      ldp    x19, x20, [x8], #16      // restore callee-saved registers
975      ldp    x21, x22, [x8], #16
976      ldp    x23, x24, [x8], #16
977      ldp    x25, x26, [x8], #16
978      ldp    x27, x28, [x8], #16
979      ldp    x29, x9, [x8], #16
980      ldr    lr, [x8]
981      mov     sp, x9
982      msr    sp_el0, x1
983      ret
984  ENDPROC(cpu_switch_to)
```

Context switch x86 assembly code

Linux

```
224  /*
225  * %eax: prev task
226  * %edx: next task
227  */
228  ENTRY(__switch_to_asm)
229      /*
230      * Save callee-saved registers
231      * This must match the order in struct inactive_task_frame
232      */
233      pushl  %ebp
234      pushl  %ebx
235      pushl  %edi
236      pushl  %esi
237
238      /* switch stack */
239      movl   %esp, TASK_threadsp(%eax)
240      movl   TASK_threadsp(%edx), %esp
241
242      /* restore callee-saved registers */
243      popl   %esi
244      popl   %edi
245      popl   %ebx
246      popl   %ebp
247
248      jmp    __switch_to
249  END(__switch_to_asm)
```

xv6

```
1  # Context switch
2  #
3  # void swtch(struct context **old, struct context *new);
4  #
5  # Save current register context in old
6  # and then load register context from new.
7
8  .globl swtch
9  swtch:
10     movl 4(%esp), %eax } Difference #1: xv6 passes
11     movl 8(%esp), %edx } parameters on stack
12
13     # Save old callee-save registers
14     pushl %ebp
15     pushl %ebx
16     pushl %esi } Difference #2: %esi & %edi
17     pushl %edi } registers are in different order
18
19     # Switch stacks
20     movl %esp, (%eax)
21     movl %edx, %esp
22
23     # Load new callee-save registers
24     popl %edi
25     popl %esi
26     popl %ebx
27     popl %ebp
28     ret
```

Process creation in Unix

- Uses a combination of *fork* and *exec* syscalls
- **Fork** creates an exact duplicate of the current process, except
 - Has a new process id
 - Parent/child processes are different
 - Return code of `fork()` command is different (...you'll see what I mean)
- **Exec** overwrites the code of the current process with that in a file
- It looks like a strange design, but it makes the command-line shell implementation clean.

Fork syscall

```
int main(int argc, char *argv[]) {  
    printf("hello world (pid:%d)\n", (int) getpid());  
    int rc = fork();  
    if (rc < 0) { // fork failed; exit  
        fprintf(stderr, "fork failed\n");  
        exit(1);  
    } else if (rc == 0) { // child (new process)  
        printf("hello, I am child (pid:%d)\n", (int) getpid());  
    } else { // parent goes down this path (main)  
        printf("hello, I am parent of %d (pid:%d)\n",  
              rc, (int) getpid());  
    }  
    return 0;  
}
```

Output:

```
hello world (pid:29146)  
hello, I am parent of 29147 (pid:29146)  
hello, I am child (pid:29147)
```

- The new (child) process continues where the parent left off.
 - It does not start from the beginning of main()
- fork returns:
 - 0 to the child process
 - the child pid to the parent
- Two processes share the same stdin, stdout, & stderr

Nondeterminism

```
int main(int argc, char *argv[]) {
    printf("hello world (pid:%d)\n", (int) getpid());
    int rc = fork();
    if (rc < 0) { // fork failed; exit
        fprintf(stderr, "fork failed\n");
        exit(1);
    } else if (rc == 0) { // child (new process)
        printf("hello, I am child (pid:%d)\n", (int) getpid());
    } else { // parent goes down this path (main)
        printf("hello, I am parent of %d (pid:%d)\n",
            rc, (int) getpid());
    }
    return 0;
}
```

Output possibility 1:

```
hello world (pid:29146)
hello, I am parent of 29147 (pid:29146)
hello, I am child (pid:29147)
```

Output possibility 2:

```
hello world (pid:29146)
hello, I am child (pid:29147)
hello, I am parent of 29147 (pid:29146)
```



- At the end of the fork syscall, the OS has two runnable processes.
- We *cannot predict* whether the OS will schedule the parent or child process to run next.
 - Depends on the runtime situation and hidden kernel implementation details.
- Thus the program's output's called *nondeterministic* or *indeterminate*.
 - Meaning it can exhibit different behavior on different runs.
- There are two output possibilities:

Nondeterminism

- Can arise when a *concurrent* program has a *race condition*, meaning:
 - Two or more things are happening at the same time,
 - It's not clear which will finish first, and
 - The output will be different depending on which finishes first.
- In the fork example, the two competing tasks were:
 - The parent process waiting to run and print
 - The child process waiting to run and print
- Race conditions can lead to difficult software bugs
 - 99% of the time it behaves one way, but sometimes it behaves another way
 - *Heisenbugs* – bugs that disappear when testing (in this case due to timing)

Can you spot the tricky bug here?

```
int main(){
    // open a file
    int fd = open(filename, O_RDWR);
    if (fd == -1) { /* Handle error */ }

    char c;
    pid_t pid = fork();
    if (pid == -1) { /* Handle error */ }
    // child
    if (pid == 0) {
        read(fd, &c, 1);
        printf("child:%c\n",c);
    }
    // parent
    else {
        read(fd, &c, 1);
        printf("parent:%c\n",c);
        do_some_work();
        // close the file
        close(fd);
    }
    return 0;
}
```

- This code is nondeterministic
- Either parent or child will print first character of file
- However, this code will also *crash* in very rare scenarios.

A *race condition* between child's read and parent's close

```
int main(){
    // open a file
    int fd = open(filename, O_RDWR);
    if (fd == -1) { /* Handle error */ }

    char c;
    pid_t pid = fork();
    if (pid == -1) { /* Handle error */ }
    // child
    if (pid == 0) {
        read(fd, &c, 1);
        printf("child:%c\n",c);
    }
    // parent
    else {
        read(fd, &c, 1);
        printf("parent:%c\n",c);
        do_some_work();
        // close the file
        close(fd);
    }
    return 0;
}
```

The child's read can happen *after* the file was closed by the parent.

- Normally, *close* will happen well after both *reads*, because *do_some_work* will be slow.
- But this is not guaranteed!

Recall that CPU exceptions are a type of interrupt

- Often caused by arithmetic errors (divide by zero), and memory violations (eg., dereferencing a null or invalid pointer)
- When **user** code causes an exception:
 - Kernel interrupt handler runs, and will likely kill the user process.
- What happens when **kernel** code causes an exception?
 - Interrupt handler will still run, but it's not clear what can be done in response.
 - On Windows, the famous “blue screen of death”
 - On Linux, a “kernel panic”
 - This is commonly seen by kernel developers, but hopefully not users.
 - This is different than the machine just freezing.
 - Kernel knows there is a problem, but doesn't know how to react.

On Windows (old & new)

Windows

An error has occurred. To continue:

Press Enter to return to Windows, or

Press CTRL+ALT+DEL to restart your computer. If you do this, you will lose any unsaved information in all open applications.

Error: 0E : 016F : BFF9B3D4

Press any key to continue _



Your PC ran into a problem and needs to restart. We're just collecting some error info, and then we'll restart for you. (0% complete)

If you'd like to know more, you can search online later for this error: HAL_INITIALIZATION_FAILED

On older Macs

You need to restart your computer. Hold down the Power button for several seconds or press the Restart button.

Veillez redémarrer votre ordinateur. Maintenez la touche de démarrage enfoncée pendant plusieurs secondes ou bien appuyez sur le bouton de réinitialisation.

Sie müssen Ihren Computer neu starten. Halten Sie dazu die Einschalttaste einige Sekunden gedrückt oder drücken Sie die Neustart-Taste.

コンピュータを再起動する必要があります。パワーボタンを数秒間押し続けるか、リセットボタンを押してください。

On Linux

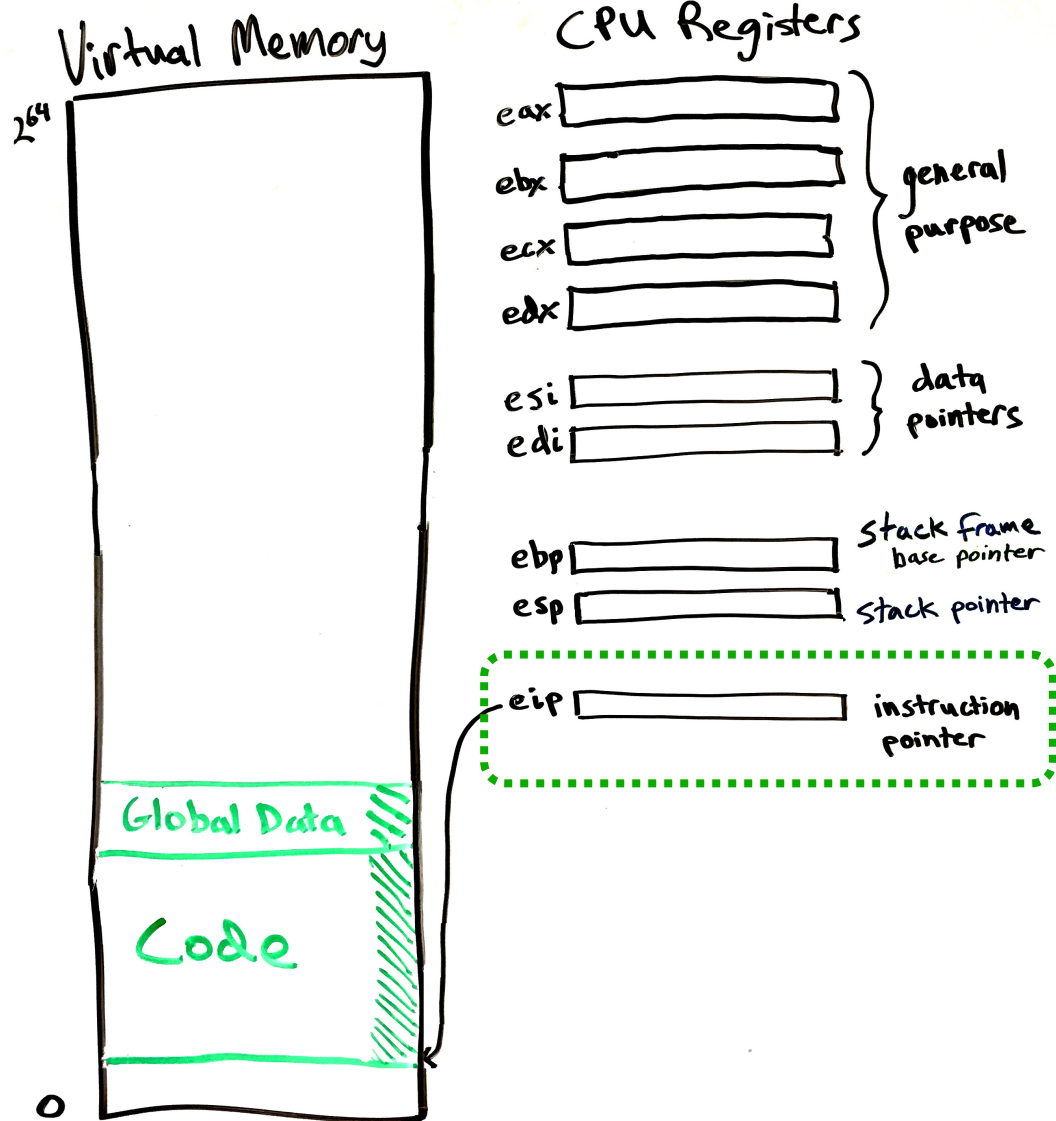
```
[<ffffffff813d059a>] ? 0xffffffff813d059a
[<ffffffff813d155b>] ? 0xffffffff813d155b
[<ffffffff813d21ea>] ? 0xffffffff813d21ea
[<ffffffff82317c2d>] ? 0xffffffff82317c2d
[<ffffffff81000296>] ? 0xffffffff81000296
[<ffffffff8104a19f>] ? 0xffffffff8104a19f
[<ffffffff822e5e81>] ? 0xffffffff822e5e81
[<ffffffff822e5647>] ? 0xffffffff822e5647
[<ffffffff81a8e853>] ? 0xffffffff81a8e853
[<ffffffff81a8e859>] ? 0xffffffff81a8e859
[<ffffffff81a9d198>] ? 0xffffffff81a9d198
[<ffffffff81a8e853>] ? 0xffffffff81a8e853
Code: 10 d6 ff ff 8b 92 00 b6 00 00 89 d2 48 8b 8f 10 d6 ff ff 8b 89 04 b6 00 00
48 c1 e1 20 48 09 d1 31 d2 49 89 c8 49 29 c0 4c 89 c0 <49> f7 f1 48 89 c8 48 85
d2 75 05 4d 39 d0 76 05 ff ce 75 be c3
RIP [<ffffffff817484bf>] 0xffffffff817484bf
RSP <ffff88014c8ffc90>
---[ end trace c384d3e911d6a1b6 ]---
Kernel panic - not syncing: Attempted to kill init! exitcode=0x0000000b

Kernel Offset: 0x0 from 0xffffffff81000000 (relocation range: 0xffffffff80000000
-0xffffffff9fffffff)
---[ end Kernel panic - not syncing: Attempted to kill init! exitcode=0x0000000b
```

Intermission recap

- xv6 OS code is written for the Intel x86 CPU architecture, but...
- Linux supports 31 different CPU architectures
 - Low-level *mechanisms* are different on each arch.
 - High-level *policies* are the same for all.
- *Fork* syscall: run once, exits twice!
- *Nondeterminism* is when a program's output is unpredictable
- OS process scheduler can create *race conditions* in programs that rely on an interaction of multiple processes.
 - These are tricky to debug, because they are sensitive to timing (*Heisenbugs*).
- *Kernel panic* occurs when OS causes an exception and can't recover

Starting a process

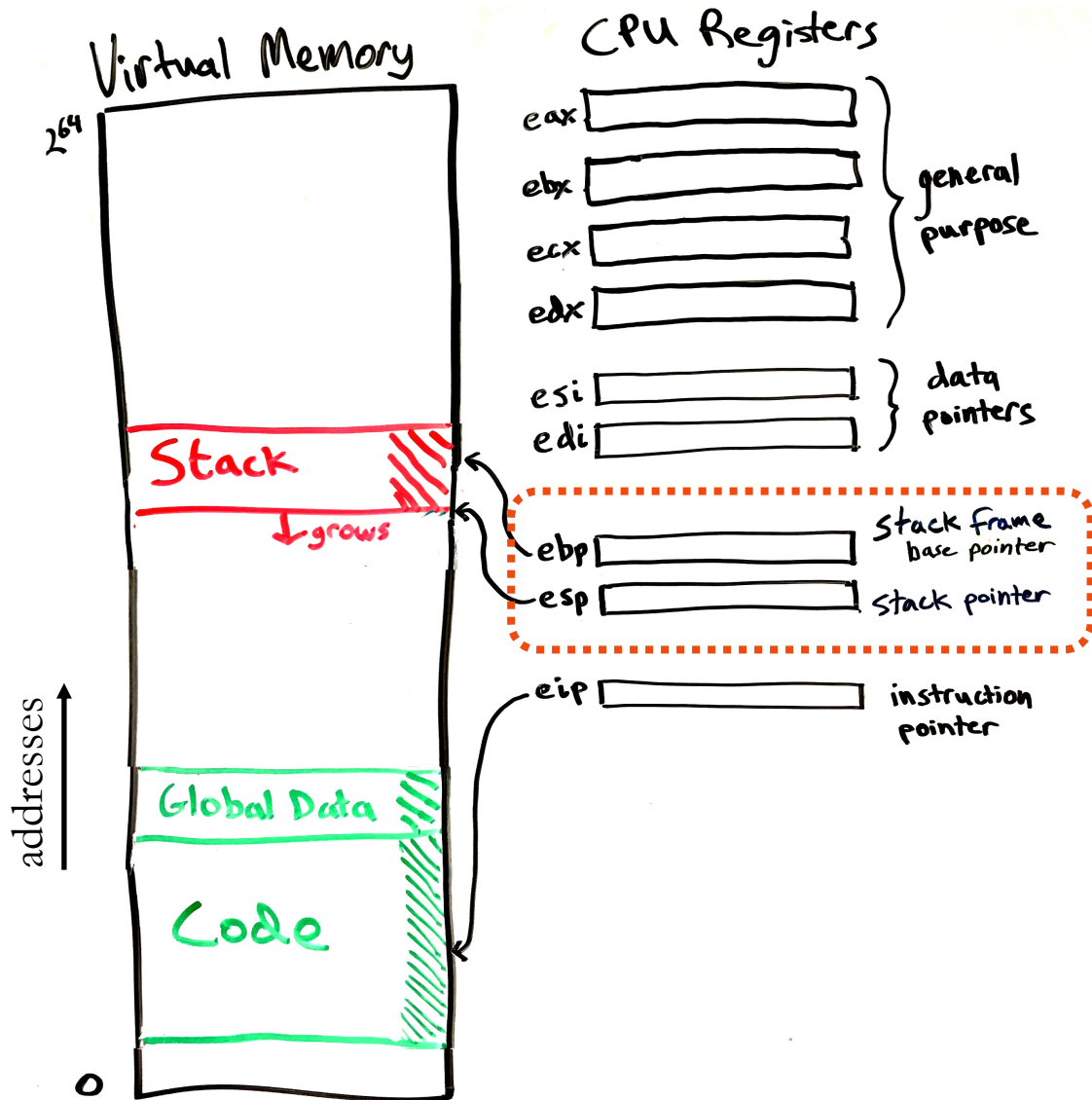


Requires just a few steps:

- Copy machine code and initial data into memory
 - In other words, copy the program's executable file into memory
- Set instruction pointer register to address of code start
 - In other words, *jump* to code start

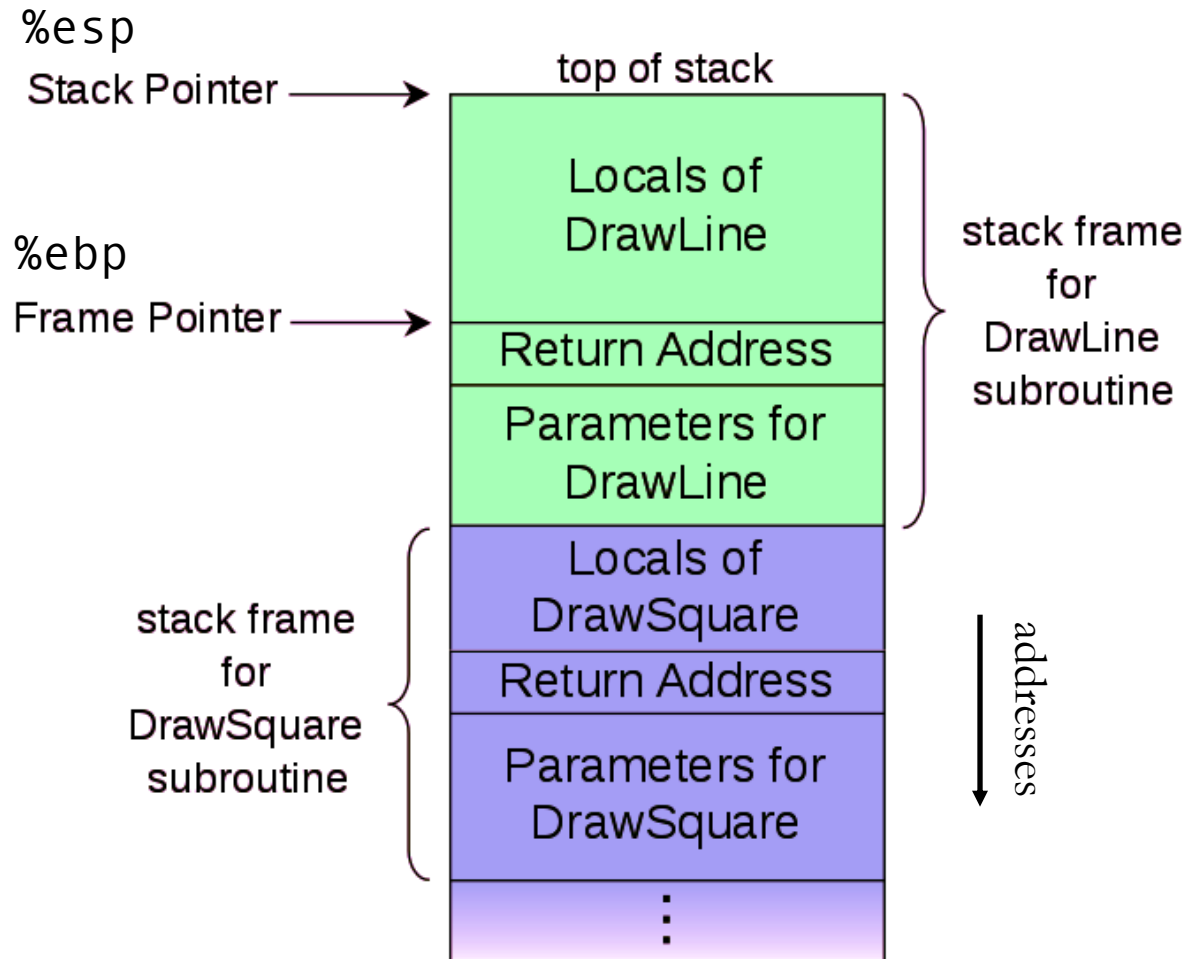
Code will use the registers and memory as necessary to perform its work.

What's this *stack* we always talk about?



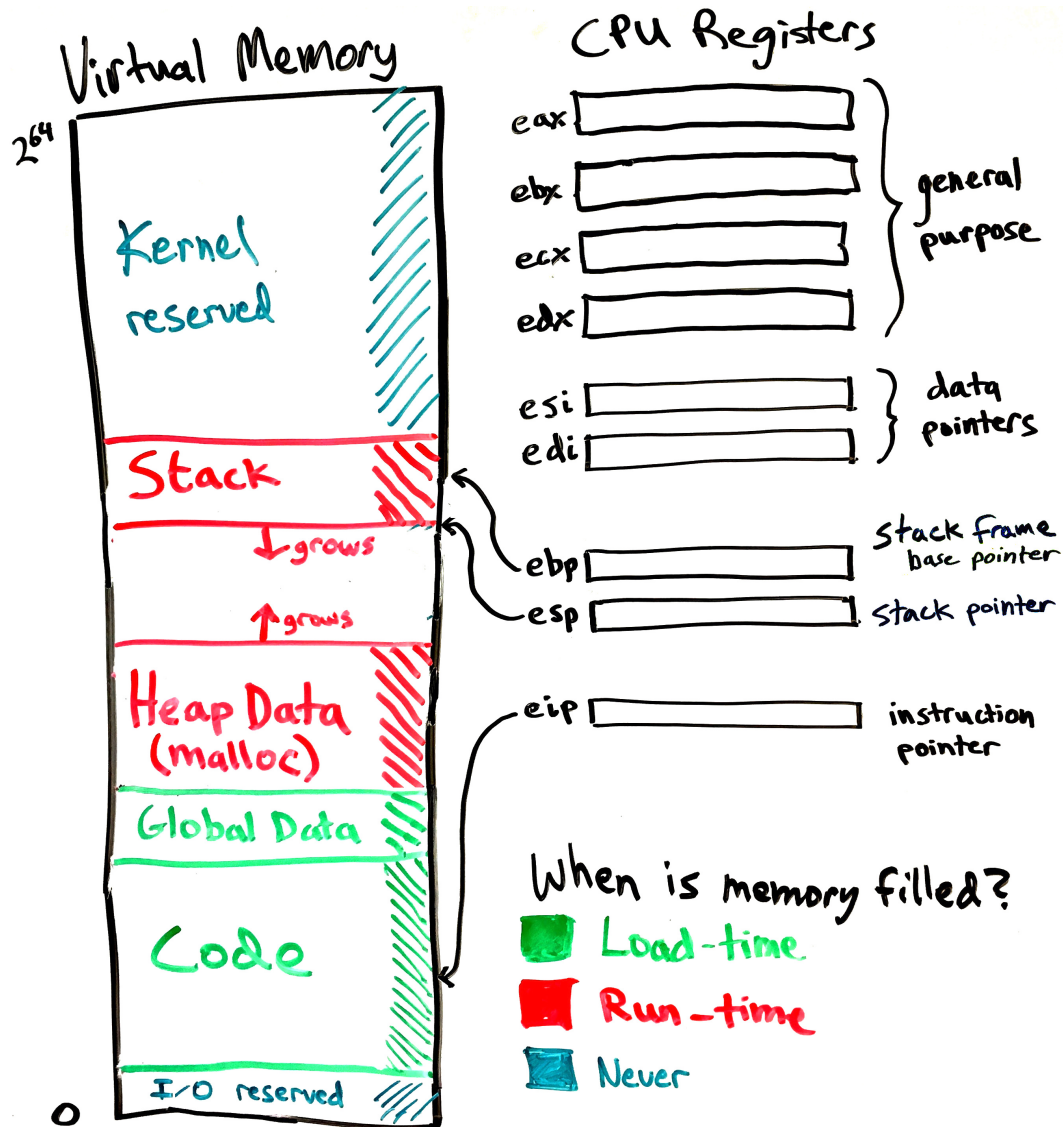
- a.k.a: execution stack, machine stack, call stack, control stack
- It's just a convenience for the assembly programmer/compiler.
 - Allows program to call subroutines and manage local variables with just a few instructions.
- Stack pointer (`%esp`) is used & automatically adjusted by:
 - `push, pop`
 - `call, ret` (return)

Using the stack for subroutines



- Greatly simplifies machine code generation for C-style functions
- Current function's local variables are on top of the stack
- To return,
 - restore caller's stack frame by restoring `%esp`, `%ebp`
 - Place function's return value in `%eax`
- DrawLine code can always **find** its *parameters* and *local variables*
 - Regardless of *when/where* the function was called, variables can be found relative to `%ebp`, the frame pointer
- In other words, the stack allows subroutines to be mutually *isolated*.

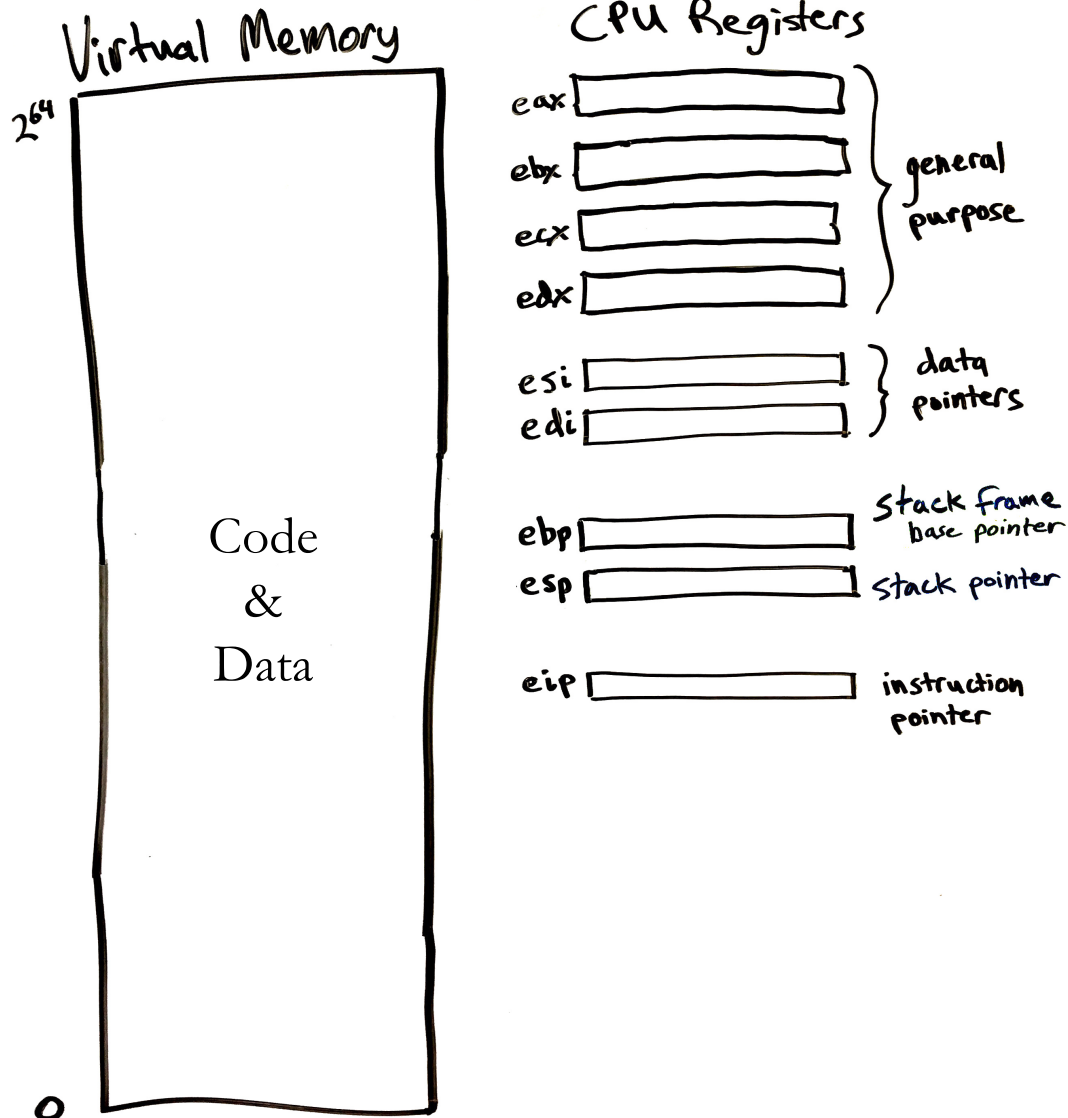
Heap memory



- Heap is just where C's *malloc* function dynamically allocates memory.
- The CPU has no notion of a special heap region.
 - Organizing memory into stack and heap is just a convention.
- Stack and Heap grow toward each other, eating free space between.

“Heap” memory has nothing to do with the “heap” self-balancing priority queue data structure.

Context Switch to change process

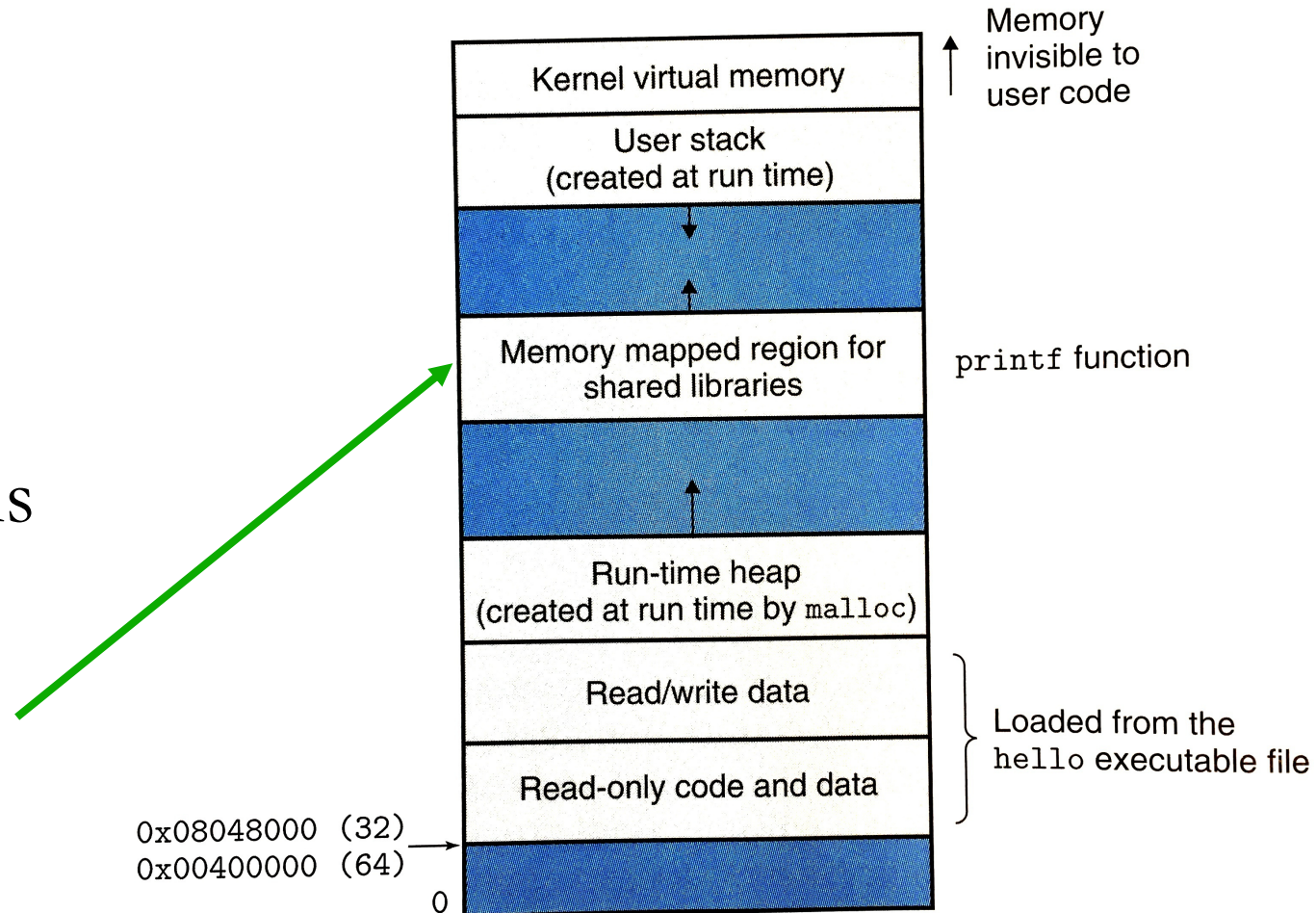


Context switch is when CPU switches from running one process to running another.

- Want context switches to be fast, to give user the illusion that processes are running simultaneously
- Need to swap out all process state
- Registers are small & fast, so they can be saved and restored
- But how to deal with memory?
 - It's big!
 - Would be too slow to copy all memory elsewhere (to disk?)

Linux process virtual memory address regions

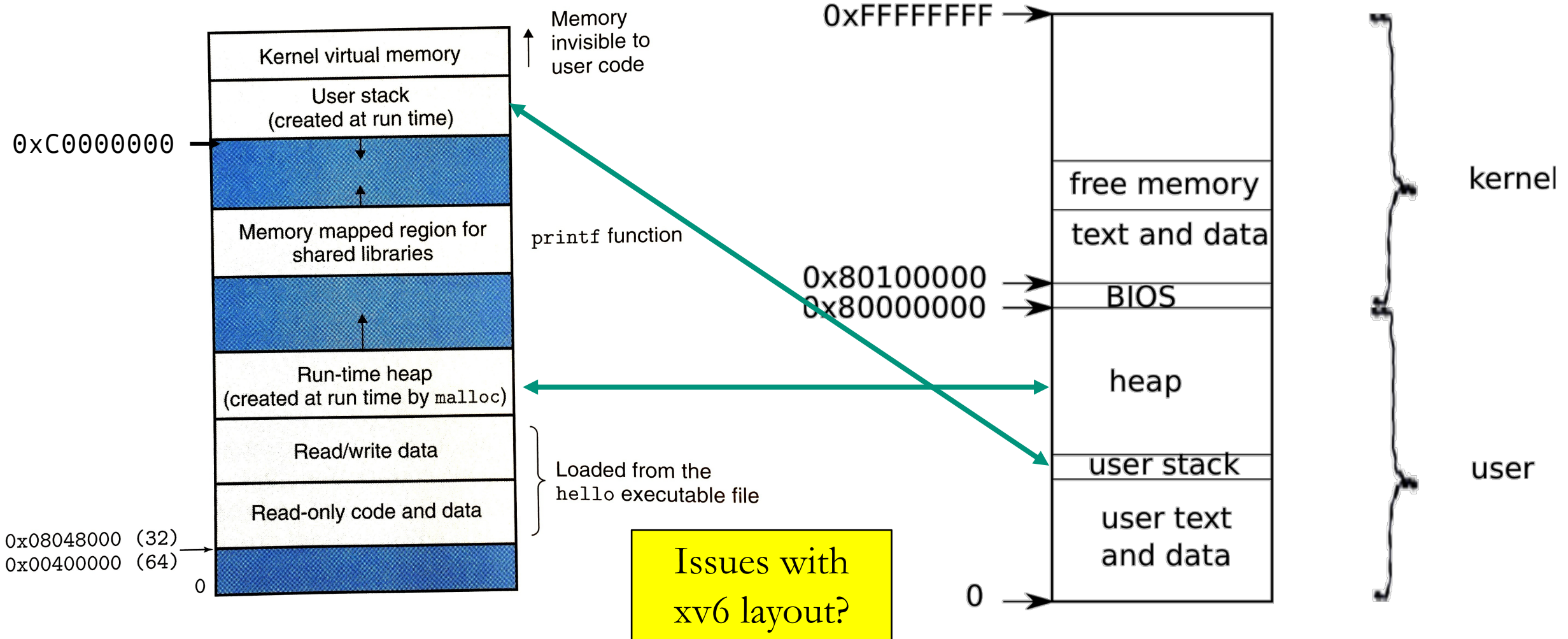
- Top of the memory range is reserved for the kernel.
 - This is actually mapped to the same physical memory for every process.
- On the PC, low memory range is reserved for I/O
- Shared libraries are not used in xv6, but they exist in modern OSes like Linux



Operating systems vary in the details

Linux process memory layout

xv6 process memory layout



Final recap

- **fork + exec** runs a program.
 - fork duplicates the current process
 - exec copies code and global data from an executable file, and creates a new empty stack.
- Stack grows from high addresses down to lower.
 - Grows larger when a function is called.
 - Shrinks when a function returns.
- Heap is a block of memory managed by C's malloc & free.

