# CS 423
# Operating System Design: The Kernel Abstraction

## Tianyin Xu

\* Thanks for Prof. Adam Bates for the slides.

Let's start with some questions.

# Overview

Process concept
- A process is the OS abstraction for executing a program with limited privileges

Dual-mode operation: user vs. kernel
- Kernel-mode: execute with complete privileges
- User-mode: execute with fewer privileges

Safe control transfer
- How do we switch from one mode to the other?

# Process Abstraction

<u>Process</u>: an instance of a program that runs with limited rights on the machine
- Thread: a sequence of instructions within a process
  - Potentially many threads per process (for now, assume 1:1)
- Address space: set of rights of a process
  - Memory that the process can access
  - Other permissions the process has (e.g., which system calls it can make, what files it can access)

**How can we permit a process to execute with only limited privileges?**

# Thought Experiment

How can we implement execution with limited privilege?

- Execute each program instruction in a simulator
- If the instruction is permitted, do the instruction
- Otherwise, stop the process
- Basic model in Javascript and other interpreted languages

How can we implement execution with limited privilege?
- Execute each program instruction in a simulator
- If the instruction is permitted, do the instruction
- Otherwise, stop the process
- Basic model in Javascript and other interpreted languages

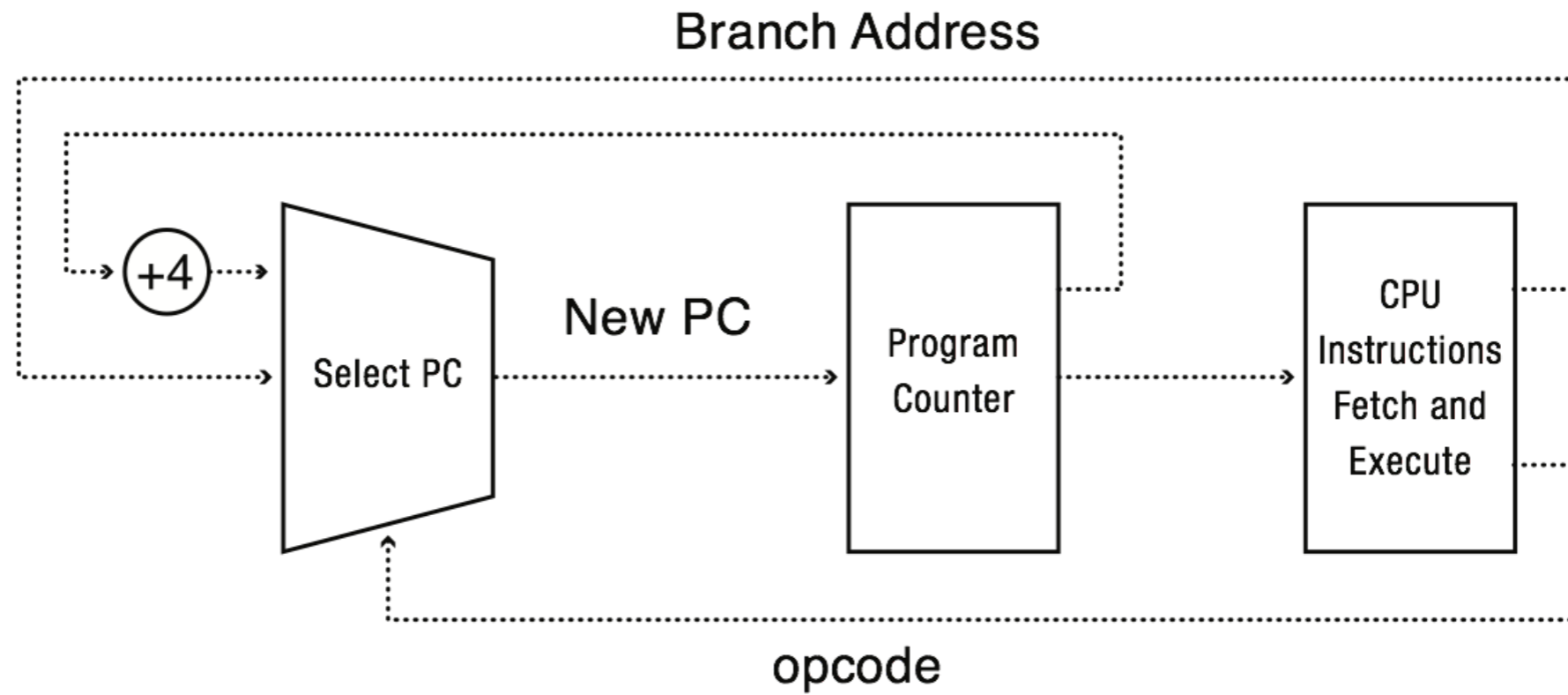**Ok… but how do we go faster?**

# Thought Experiment

How can we implement execution with limited privilege?

- Execute each program instruction in a simulator
- If the instruction is permitted, do the instruction
- Otherwise, stop the process
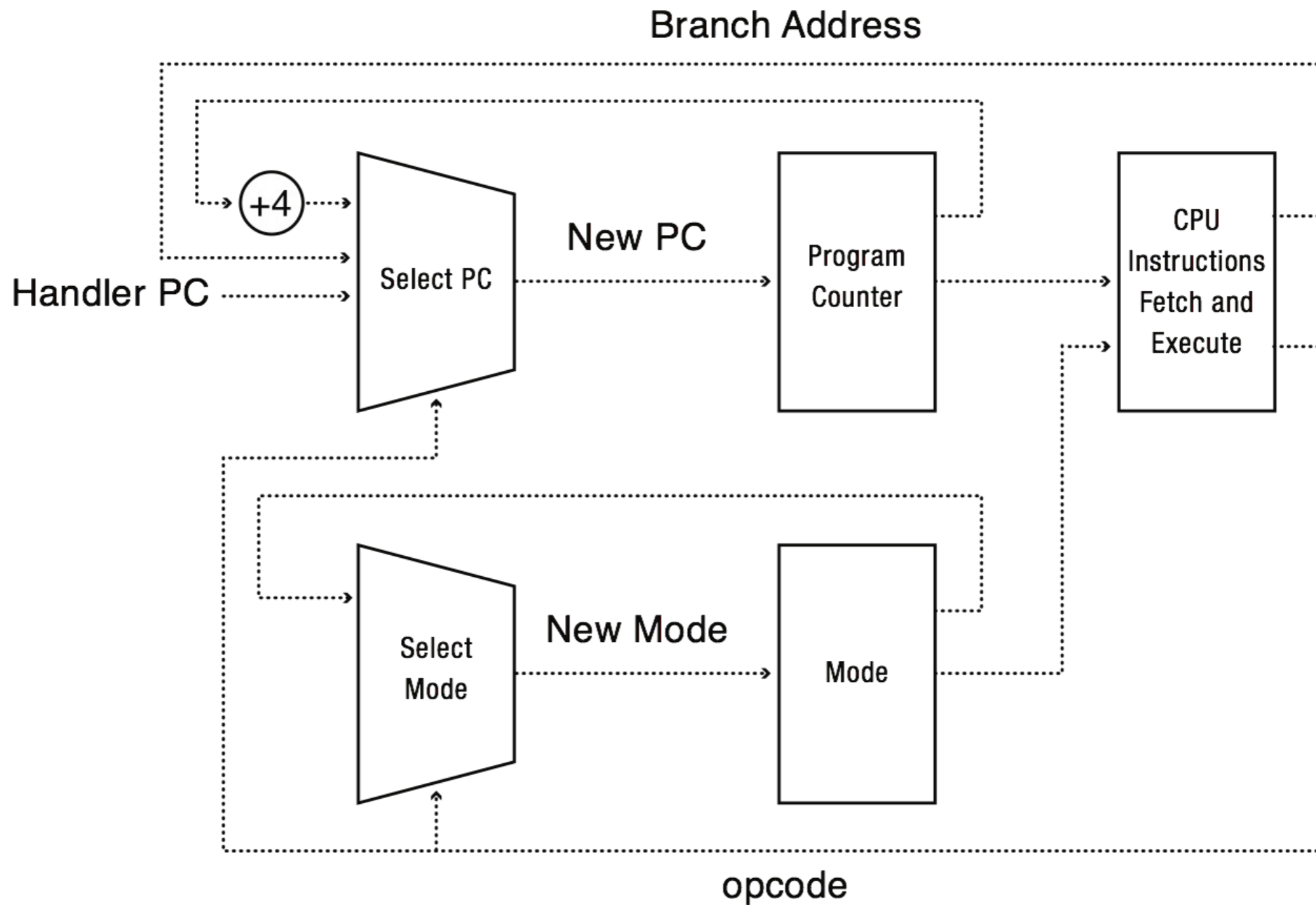- Basic model in Javascript and other interpreted languages

**Ok… but how do we go faster?**

- Run the unprivileged code directly on the CPU!

Branch Address

New PC

+4

Select PC

Program Counter

CPU Instructions Fetch and Execute

opcode

Privileged instructions
- Available to kernel
- Not available to user code

Limits on memory accesses
- To prevent user code from overwriting the kernel

Timer
- To regain control from a user program in a loop

Safe way to switch from user mode to kernel mode, and vice versa

Examples?

What should happen if a user program attempts to execute a privileged instruction?

# User->Kernel Switches

How/when do we switch from user to kernel mode?

1. Interrupts
   - Triggered by timer and I/O devices
2. Exceptions
   - Triggered by unexpected program behavior
   - Or malicious behavior!
3. System calls (aka protected procedure call)
   - Request by program for kernel to do some operation on its behalf
   - Only limited # of very carefully coded entry points

**How does the OS know when a process is in an infinite loop?**

# Hardware Timer

Hardware device that periodically interrupts the processor

- Returns control to the kernel handler
- Interrupt frequency set by the kernel
  Not by user code!
- Interrupts can be temporarily deferred
  Not by user code! Interrupt deferral crucial for implementing mutual exclusion

# Kernel->User Switches

How/when do we switch from kernel to user mode?

1. New process/new thread start
   - Jump to first instruction in program/thread
2. Return from interrupt, exception, system call
   - Resume suspended execution (return to PC)
3. Process/thread context switch
   - Resume some other process (return to PC)
4. User-level upcall (UNIX signal)
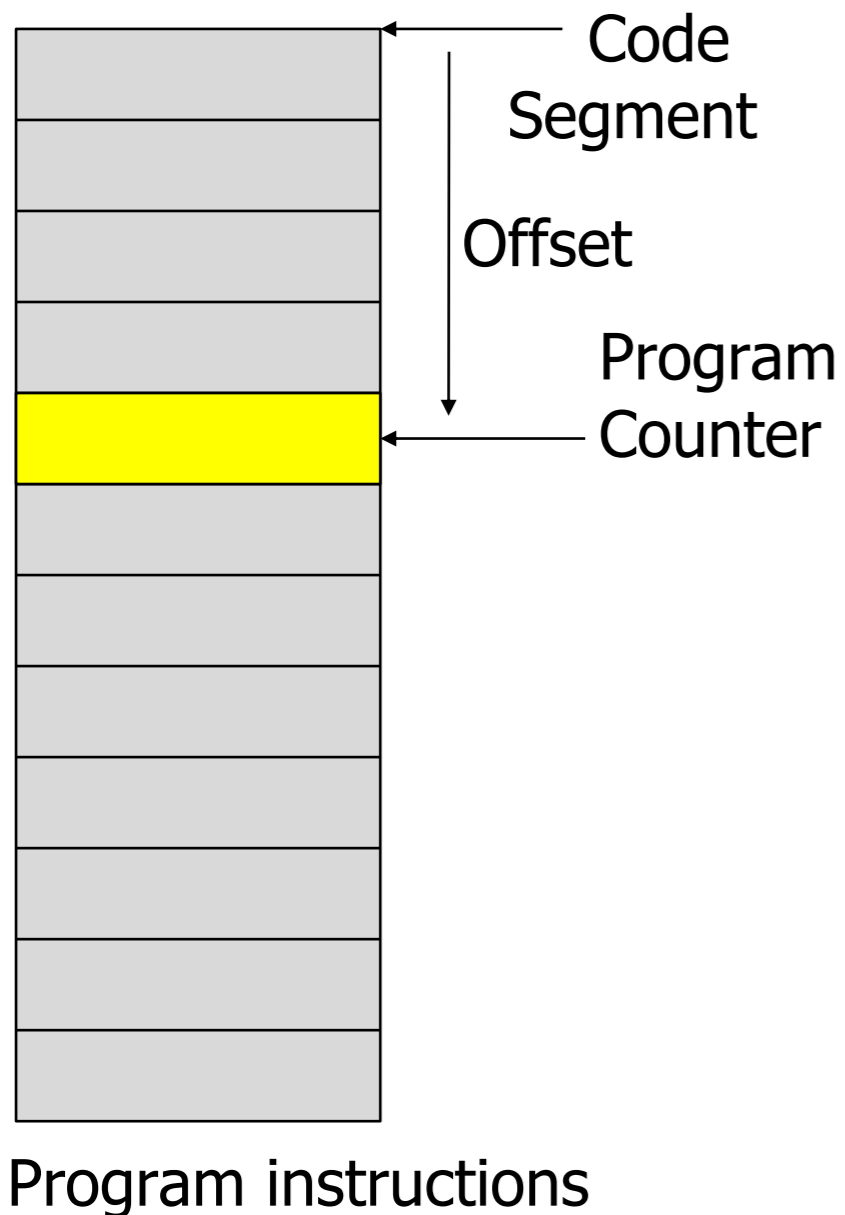   - Asynchronous notification to user program

# CPU State

What is the CPU's behavior defined by at any given moment?
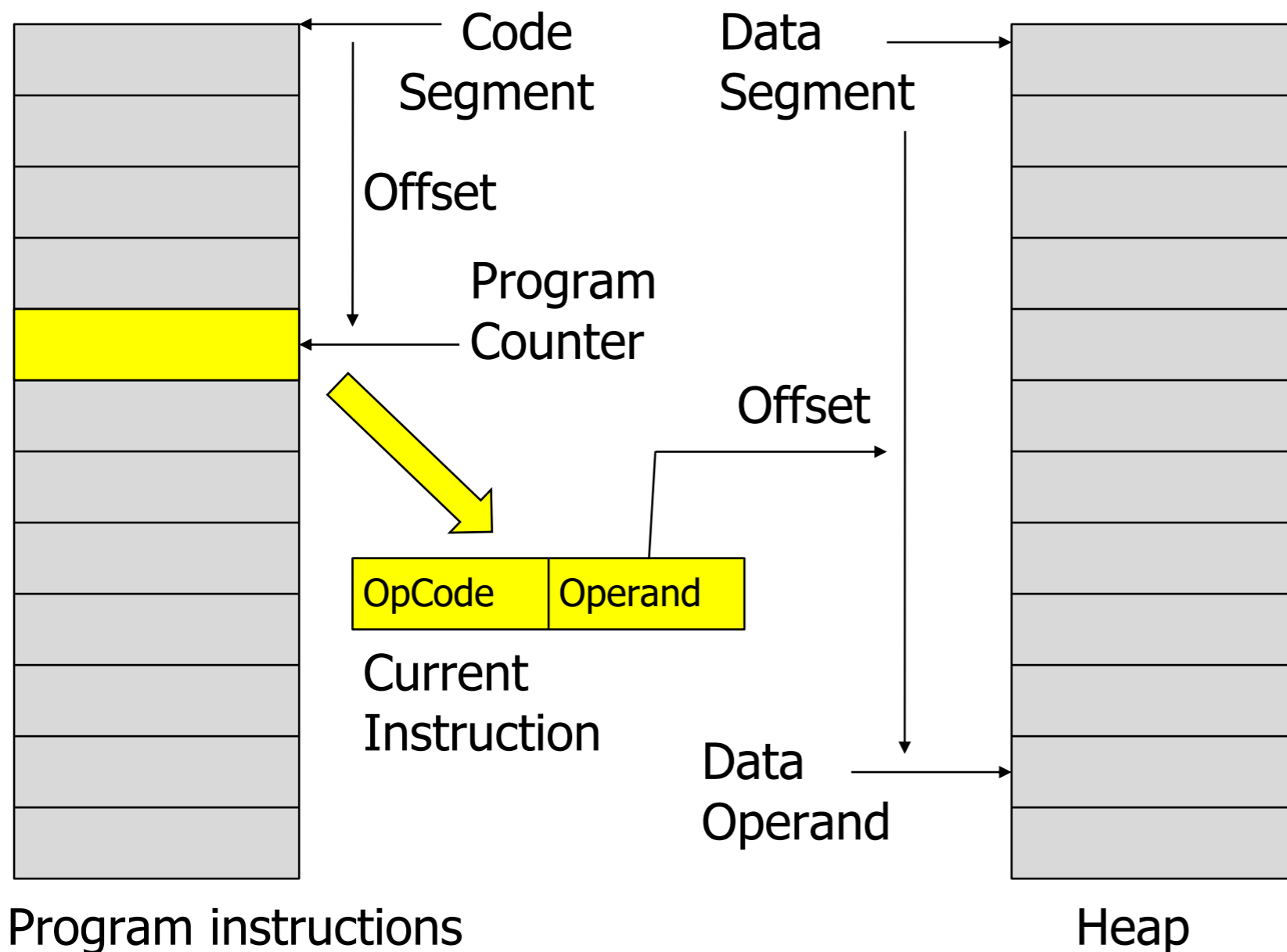
# CPU State

What is the CPU's behavior defined by at any given moment?

Code Segment

Offset

Program Counter

Program instructions

What is the CPU's behavior defined by at any given moment?



Program instructions

Heap

# CPU State

## What is the CPU's behavior defined by at any given moment?

Code Segment

Offset

Program Counter

Current Instruction

OpCode | Operand

Data Segment

Offset

Data Operand

Stack Segment

Offset

Stack Pointer

Program instructions
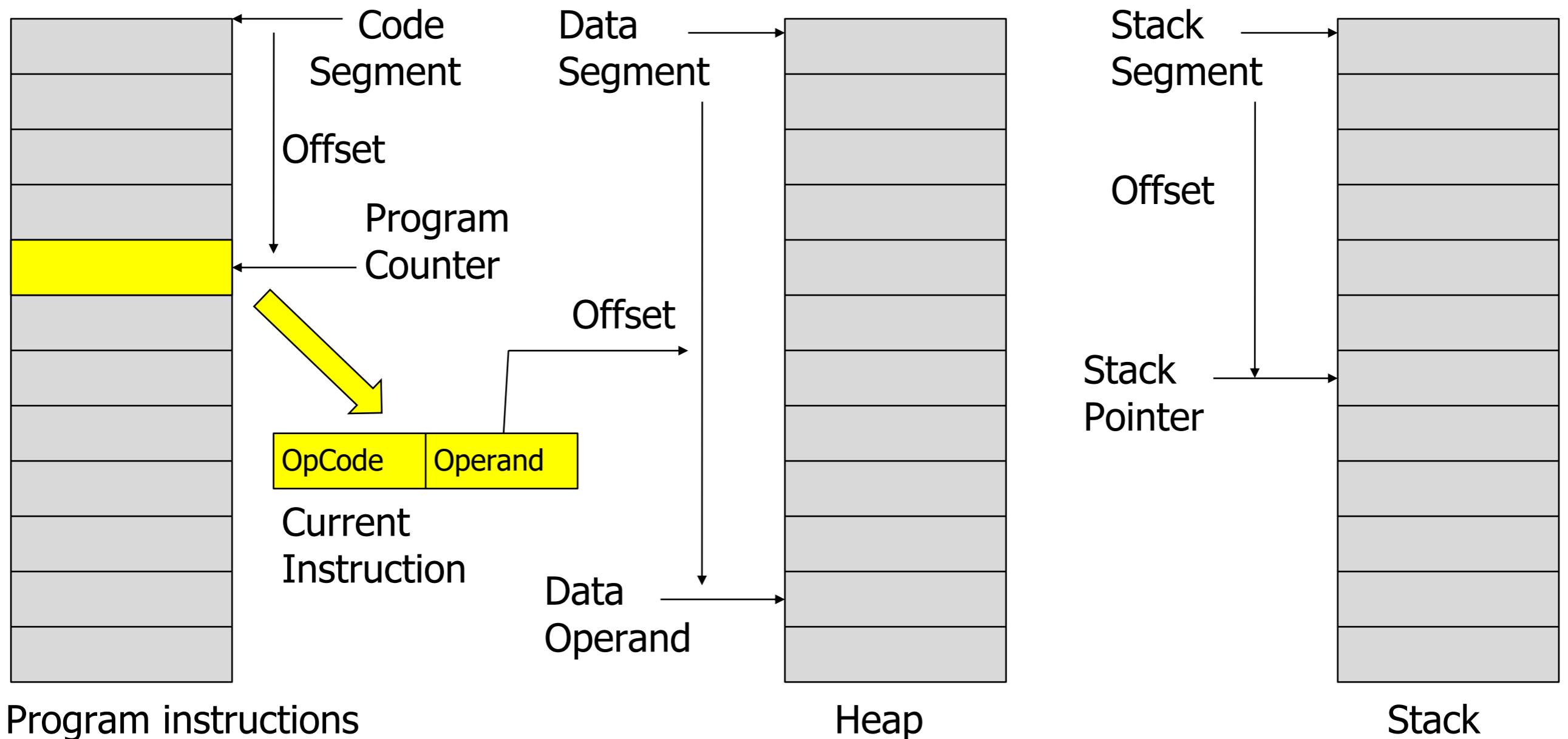
Heap

Stack

# CPU State

What is the CPU's behavior defined by at any given moment?

Registers

Code Segment

Offset

Program Counter

| OpCode | Operand |
|--------|---------|

Current Instruction

Program instructions

Data Segment

Offset

Data Operand

Heap

Stack Segment

Offset

Stack Pointer

Stack

**What defines the STATE of the CPU?**

Registers

Code Segment

Offset

Program Counter

| OpCode | Operand |

Current Instruction

Data Segment

Offset

Data Operand

Stack Segment

Offset

Stack Pointer

Program instructions

Heap

Stack

# What's a 'real' CPU?

**What's the STATE of a real CPU?**

**Registers**

**Code Segment**

Offset

**Program Counter**

Current Instruction

| OpCode | Operand |

**Data Segment**

Offset

Offset

Data Operand

Program instructions

Heap

**Stack Segment**

Offset

**Stack Pointer**

Stack

# The Context Switch

# Process Control Block

The state for processes that are not running on the CPU are maintained in the Process Control Block (PCB) data structure

**Process Table**

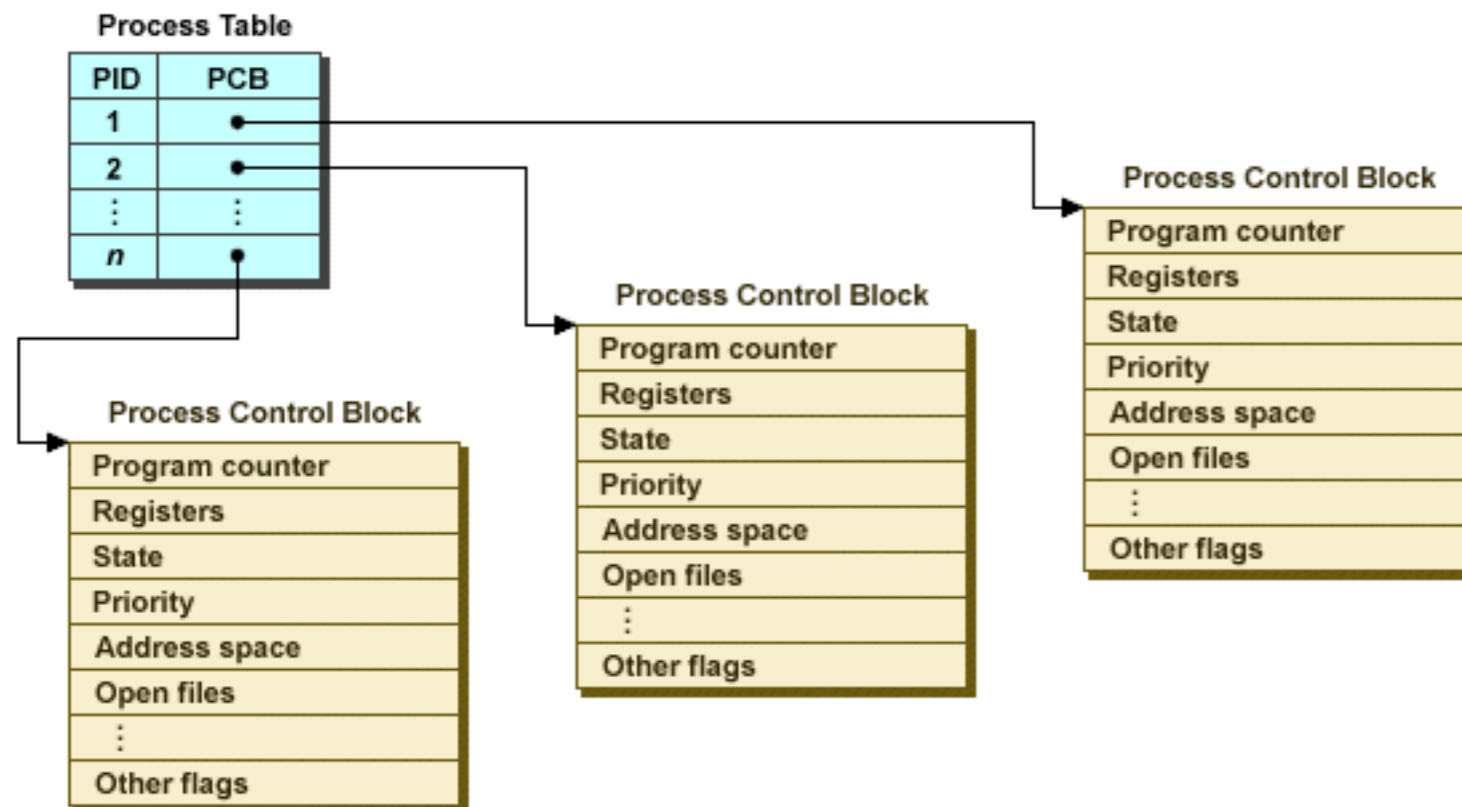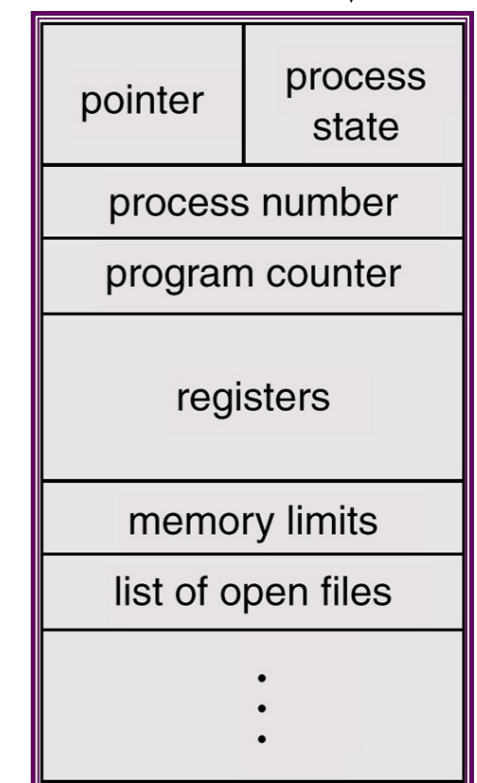| PID | PCB |
|-----|-----|
| 1   | •   |
| 2   | •   |
| ⋮   | ⋮   |
| n   | •   |

**Process Control Block**

- Program counter
- Registers
- State
- Priority
- Address space
- Open files
- ⋮
- Other flags

**Process Control Block**

- Program counter
- Registers
- State
- Priority
- Address space
- Open files
- ⋮
- Other flags

**Process Control Block**

- Program counter
- Registers
- State
- Priority
- Address space
- Open files
- ⋮
- Other flags

Updated during context switch

| pointer | process state |
|---------|---------------|
| process number | |
| program counter | |
| registers | |
| memory limits | |
| list of open files | |
| ⋮ | |

*An alternate PCB diagram*

# The Context Switch

# The Context Switch

**Registers**

**Code Segment**

Offset

**Program Counter**

OpCode    Operand

Program instructions

**Stack Segment**

**Stack Pointer**

Stack

**Data Segment**

Data Operand

Heap

Note: In **thread** context switches, heap is not switched!

**Load State (Context)**

**Save State (Context)**

**Registers**

**Code Segment**

Offset

**Program Counter**

OpCode    Operand

Program instructions

**Stack Segment**

**Stack Pointer**

Stack

# The Context Switch

# Thread Context Switch

Note: In **thread** context switches, heap is not switched!

**Registers**

**Data Segment**

**Code Segment**

Offset

**Program Counter**

OpCode    Operand

**Stack Segment**

**Stack Pointer**

Data

Program instructions

**Global Variables**

**Load State (Context)**

## So who does the context switch, and when???

**Save State (Context)**

**Registers**

**Code Segment**

fset

**Program Counter**

OpCode    Operand

**Stack Segment**

**Stack Pointer**

**Local Variables**

Program instructions

Stack

# Thread Context Switch

**Registers**

**Code Segment**

Offset

**Program Counter**

OpCode  Operand

Program instructions

**Stack Segment**
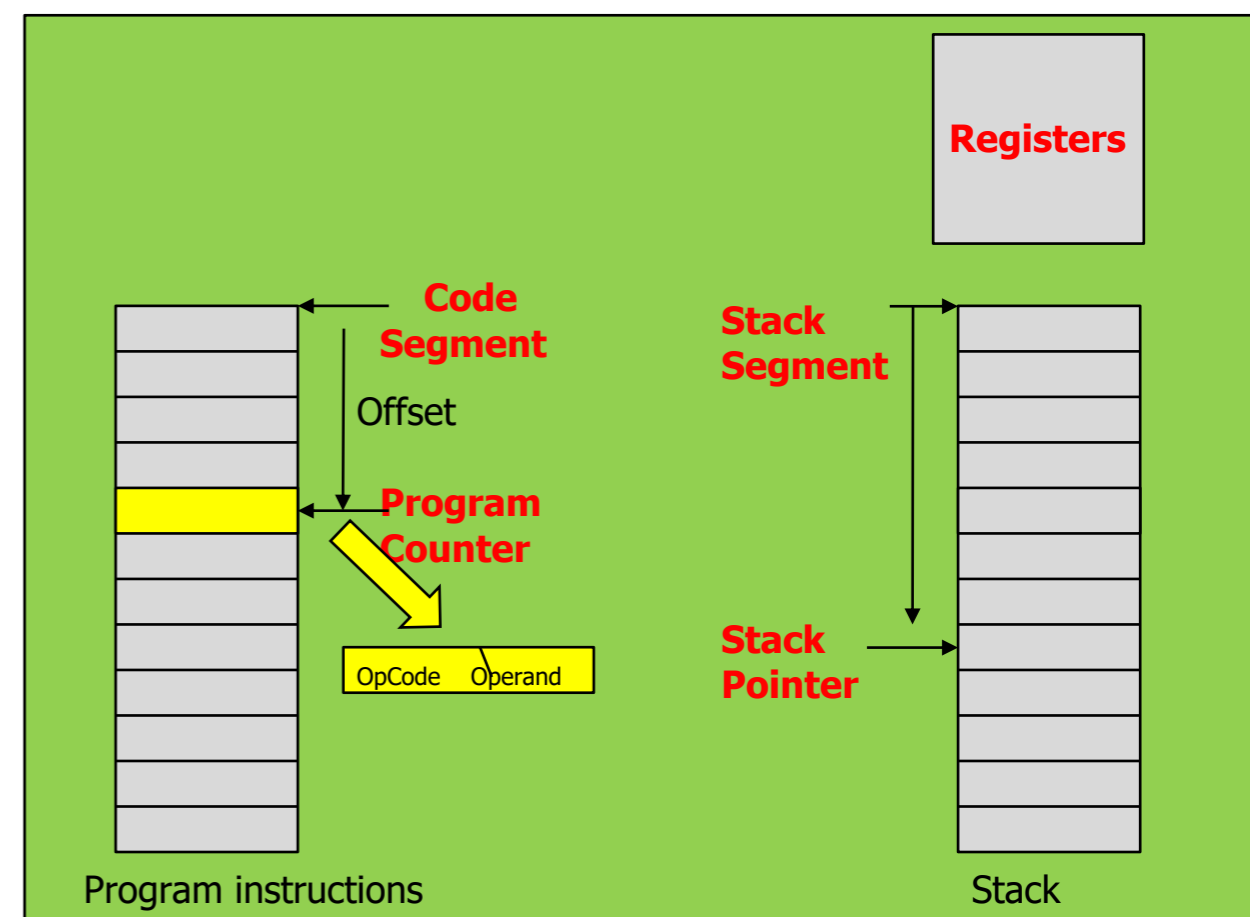
**Stack Pointer**

**Data Segment**

**Global Variables**
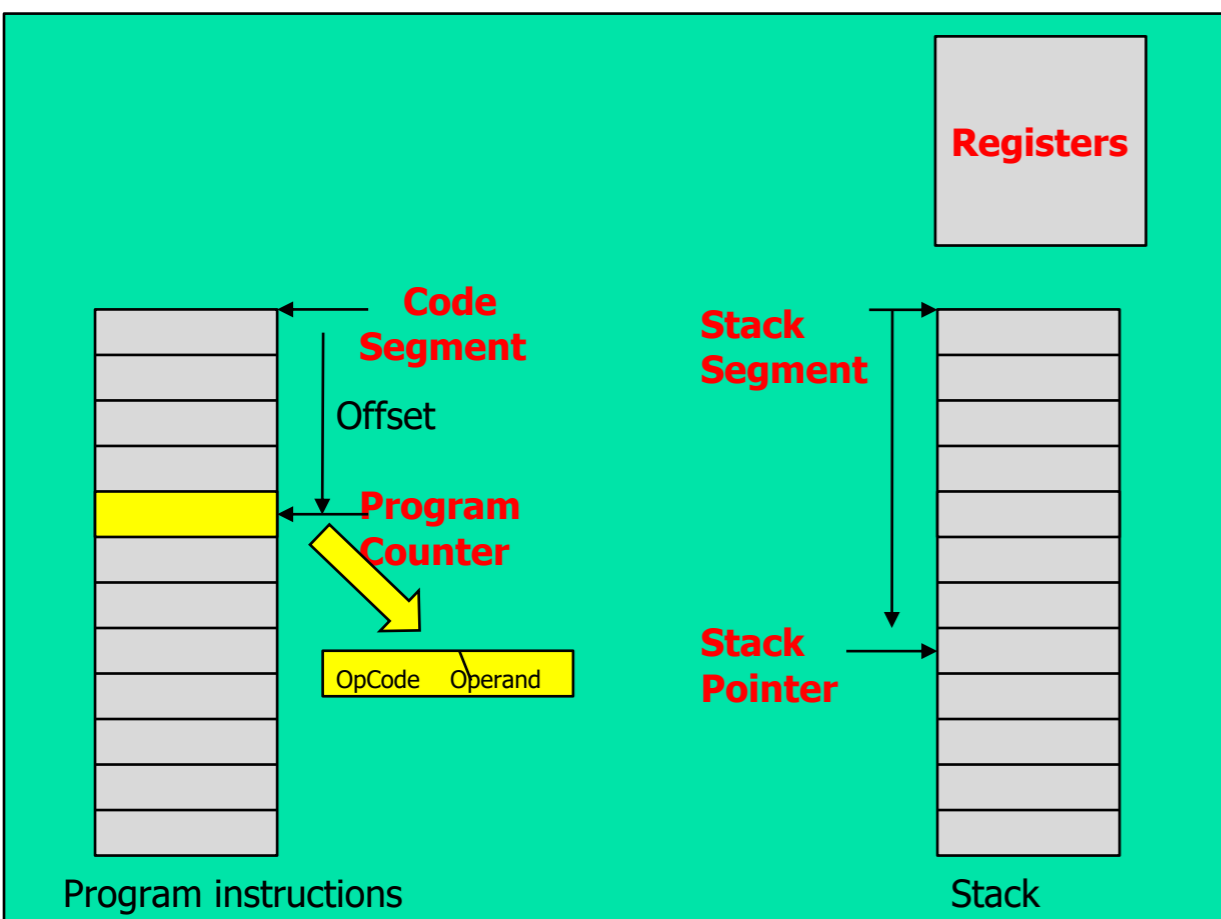
Data Operand

Heap

Note: In **thread** context switches, heap is not switched!

**Load State (Context)**

**Save State (Context)**

**Solution 1: An Interrupt**

**Registers**

**Code Segment**

Offset

**Program Counter**

OpCode  Operand

Program instructions

**Stack Segment**

**Local Variables**

**Stack Pointer**

Stack

# CTX Switch: Interrupt

**Running Thread**

Registers

Code Segment

Offset

Program Counter

Stack Segment

Stack Pointer

Program instructions

Stack

Registers

Code Segment

Offset

Program Counter

Stack Segment

Stack Pointer

Program instructions

Stack

# CTX Switch: Interrupt

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Interrupt**

Save PC on thread stack
Jump to Interrupt handler

# CTX Switch: Interrupt



**Registers**

Code Segment

Offset

Program Counter

Stack Segment

Stack Pointer

Program instructions

Stack

**Registers**

Code Segment

Offset

Program Counter

Stack Segment

Stack Pointer

Program instructions

Stack

Save PC on thread stack
Jump to Interrupt handler

Handler
- Save thread state in thread control block
  (SP, registers, segment pointers, …)

Thread Control Block

# CTX Switch: Interrupt



Registers

Code Segment
Offset
Program Counter

Stack Segment
Stack Pointer

Program instructions

Stack

Registers

Code Segment
Offset
Program Counter

Stack Segment
Stack Pointer

Program instructions

Stack

Save PC on thread stack
Jump to Interrupt handler

Handler
- Save thread state in thread control block
  (SP, registers, segment pointers, …)
- Choose next thread
- Load thread state from control block

Thread Control Block

Thread Control Block

# CTX Switch: Interrupt

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

Save PC on thread stack
Jump to Interrupt handler

Thread Control Block

Handler
- Save thread state in thread control block
  (SP, registers, segment pointers, …)
- Choose next thread
- Load thread state from control block
- Pop PC from thread stack (return from handler)

Thread Control Block

# CTX Switch: Interrupt



**Registers**

Code Segment

Offset

Program Counter

Stack Segment

Stack Pointer

Program instructions

Stack

**Registers**

Code Segment

Offset

Program Counter

Stack Segment

Stack Pointer

Program instructions

Stack

Save PC on thread stack
Jump to Interrupt handler

Handler
- Save thread state in thread control block
  (SP, registers, segment pointers, …)
- Choose next thread
- Load thread state from control block
- Pop PC from thread stack (return from handler)

Thread Control Block

Thread Control Block

**Where does it return?**

# CTX Switch: Interrupt

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

Save PC on thread stack
Jump to Interrupt handler

Thread Control Block

Handler
- Save thread state in thread control block
  (SP, registers, segment pointers, …)
- Choose next thread
- Load thread state from control block
- Pop PC from thread stack (return from handler)

Thread Control Block

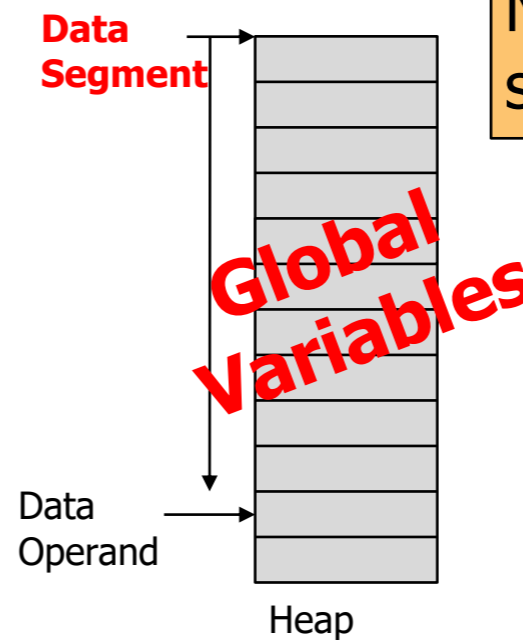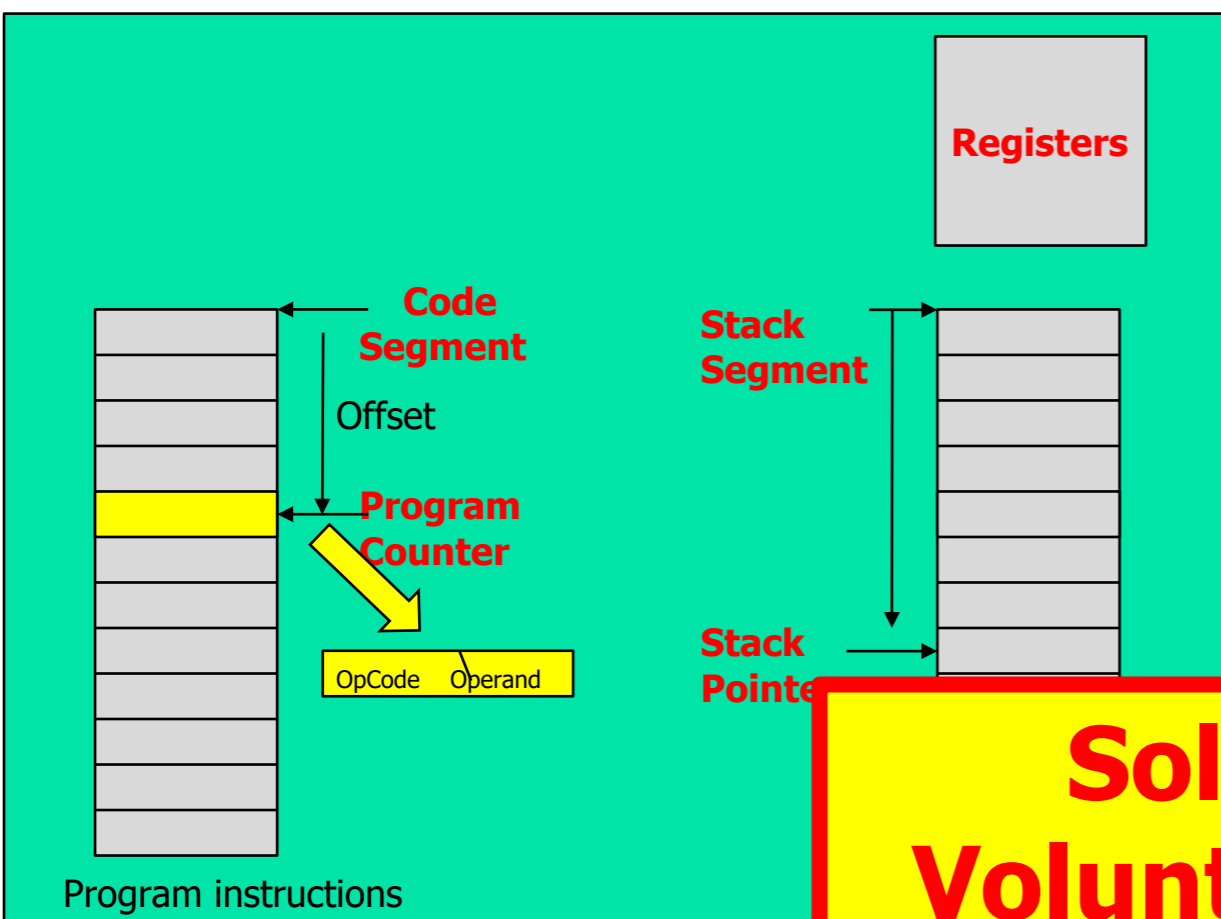**Where does it return?**

What are some examples of context switches due to interrupts?

- **Clock Interrupt:** Task exceeds its time slice

- **I/O Interrupt**: Waiting processes may be preempted

- **Memory Fault:** CPU attempts to access a virtual memory address that is not in main memory. OS may resume execution of another process while retrieving the block, then moves process to ready state.

# Thread Context Switch

**Data Segment**

**Registers**

**Global Variables**

**Load State (Context)**

**Code Segment**

Offset

**Program Counter**

OpCode    Operand

**Stack Segment**

**Stack Pointer**

Data Operand

Heap

Program instructions

**Solution 2: Voluntary yield()**

**Save State (Context)**

**Registers**

**Code Segment**

Offset

**Program Counter**

OpCode    Operand

**Stack Segment**

**Local Variables**

**Stack Pointer**

Program instructions                                   Stack

**Running Thread**

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

# CTX Switch: Yield

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**yield()**

Save PC on thread stack
Jump to yield() function

# CTX Switch: Yield



Registers

Code Segment
Offset
Program Counter

Stack Segment

Stack Pointer

Program instructions

Stack

Registers

Code Segment
Offset
Program Counter

Stack Segment

Stack Pointer

Program instructions

Stack

Save PC on thread stack
Jump to yield() function

yield()
- Save thread state in thread control block
  (SP, registers, segment pointers, …)

Thread Control Block

# CTX Switch: Yield



Save PC on thread stack
Jump to yield() function

yield()
- Save thread state in thread control block
  (SP, registers, segment pointers, …)
- Choose next thread

Thread
Control
Block

# CTX Switch: Yield



**Registers**

**Code Segment**
Offset
**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Registers**

**Code Segment**
Offset
**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

Save PC on thread stack
Jump to yield() function

yield()
- Save thread state in thread control block
  (SP, registers, segment pointers, ...)
- Choose next thread
- Load thread state from control block

Thread Control Block

Thread Control Block

# CTX Switch: Yield



**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Registers**

**Code Segment**

Offset

**Program Counter**

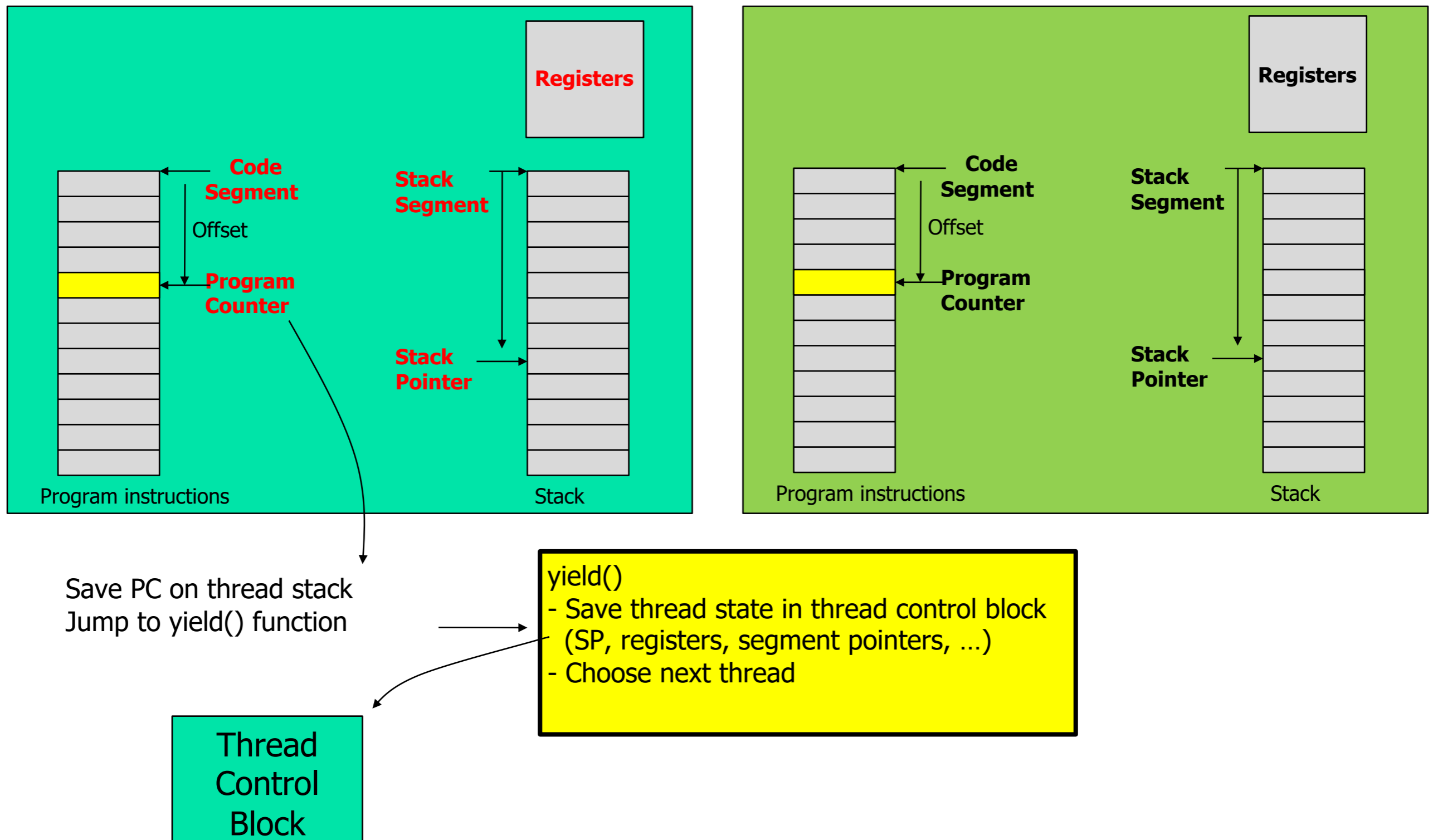**Stack Segment**

**Stack Pointer**

Program instructions

Stack

Save PC on thread stack
Jump to yield() function

yield()
- Save thread state in thread control block
  (SP, registers, segment pointers, …)
- Choose next thread
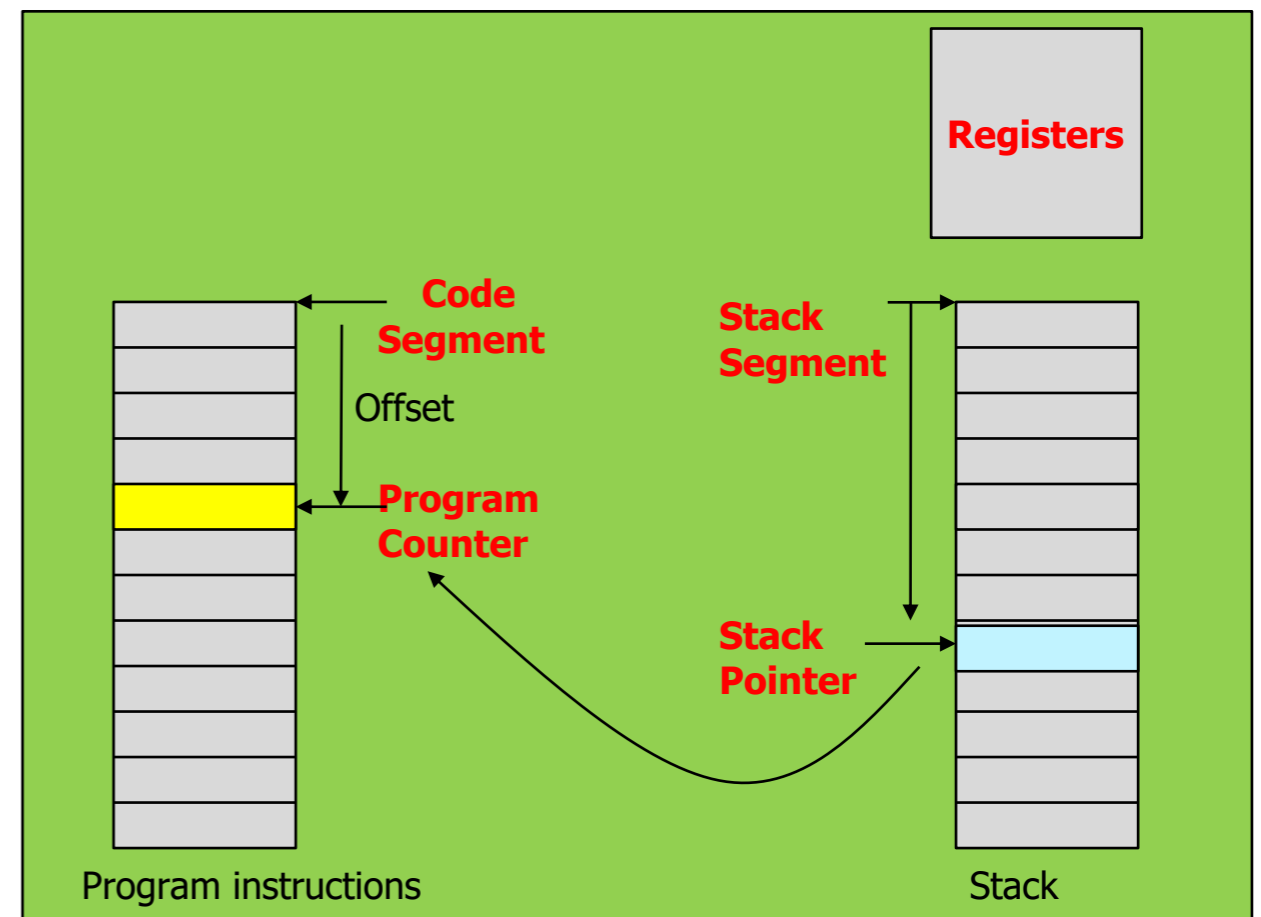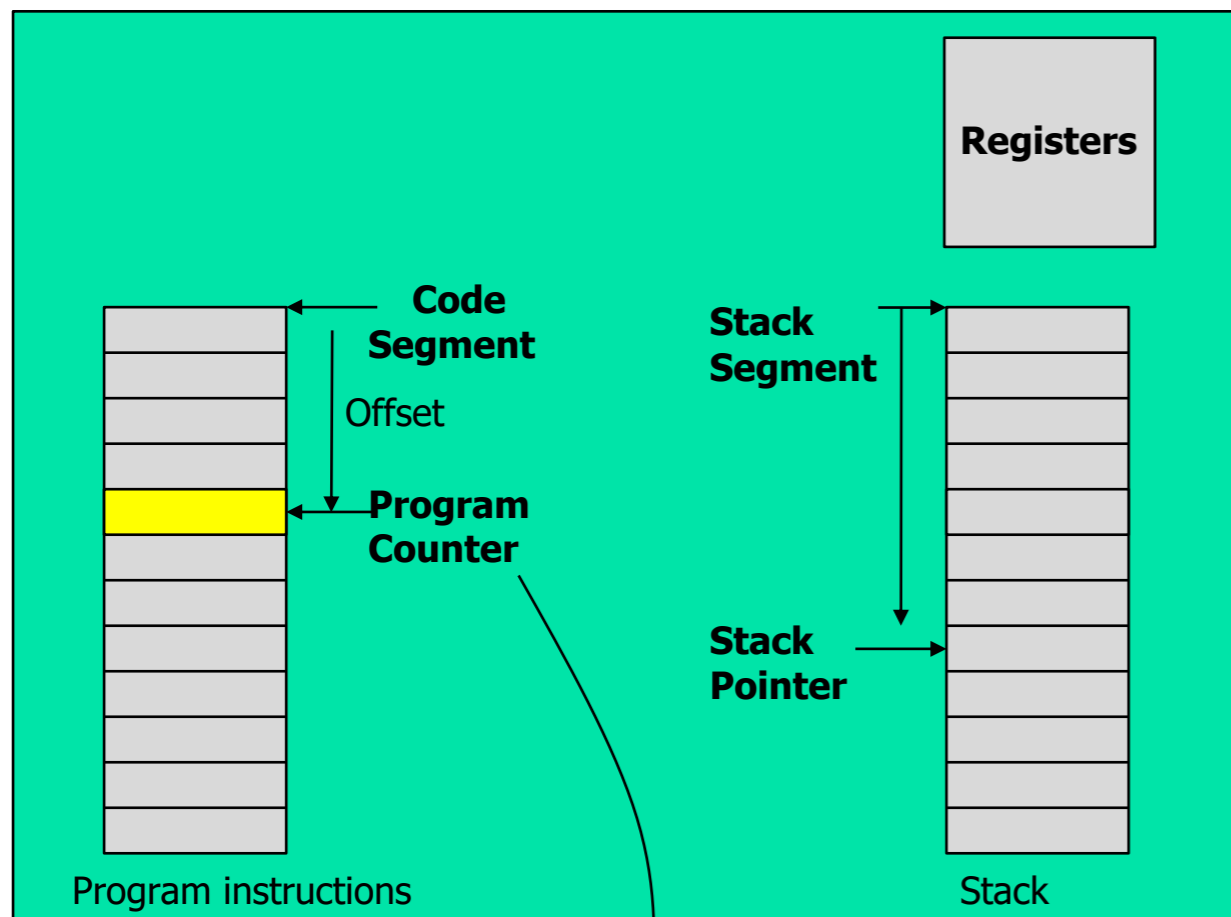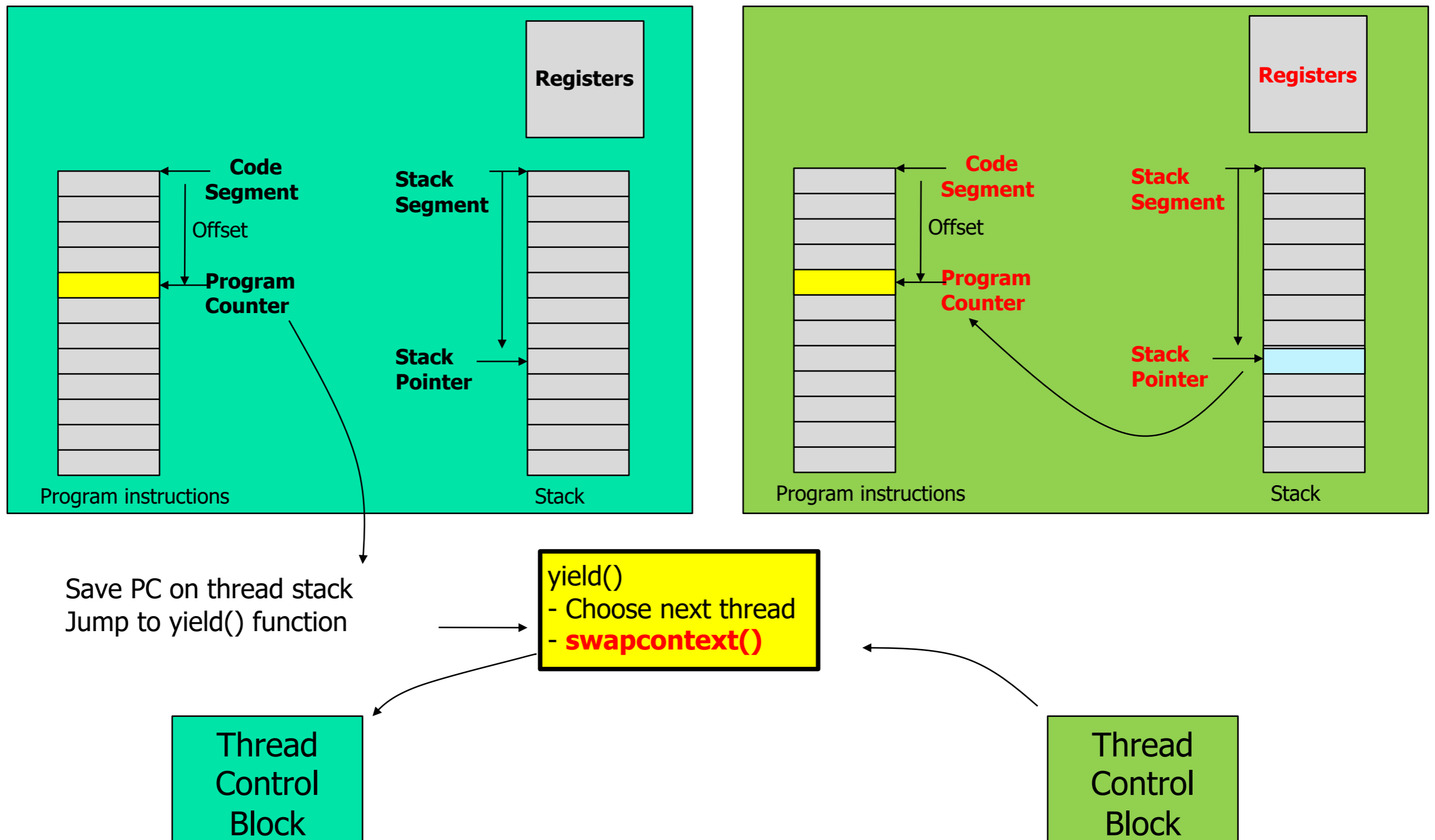- Load thread state from control block
- Pop PC from thread stack (return from handler)

Thread Control Block

Thread Control Block

# CTX Switch: Yield



Registers

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

Save PC on thread stack
Jump to yield() function

yield()
- Choose next thread
- **swapcontext()**

Thread Control Block

Thread Control Block

# Scheduler

# Scheduler

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Point**

Program instructions

Stack

**Registers**

**Code Segment**

Offset

**Program Counter**

**Stack Segment**

**Stack Pointer**

Program instructions

Stack

**Where is the Scheduling Policy?**

Save PC on thread stack
Jump to yield() function

yield()
- **NextThreadID = scheduler()**
- swapcontext()

Maintains a sorted queue of ready threads

Thread Control Block

Thread Control Block