

# Mechanisms for entering the system

Yolanda Becerra Fontal  
Juan José Costa Prats

Facultat d'Informàtica de Barcelona (FIB)  
Universitat Politècnica de Catalunya (UPC)  
BarcelonaTech  
2015-2016 QT

# Content

- Introduction
- Mechanisms for entering the system
  - Initialization
  - Management
  - Example
- Procedure for entering the system
- Procedure to exit from system
- Exceptions
- Interrupts
- System calls
- Summary

# Introduction

- OS implements access to machine resources
  - Isolates users from low-level machine-dependent code
  - Groups common code for all users: save disk space
  - Implements resource allocation policies
    - Arbitrates the usage of the machine resources in multi-user and multiprogrammed environments
  - Prevents machine and other users from user damage
    - Some instructions can not be executed by user codes: I/O instructions, halt,...

# Privilege levels (I)

- Requirement:
  - Prevent users from direct access to resources
    - Ask the OS for services
- Privilege instructions
  - Instructions that only can execute the OS
  - HW support is needed
  - When a privilege instruction is executed, the hw checks if it is executing system code
    - If not → exception
- How to distinguish user code from system code?
  - Privilege levels
    - At least 2 different levels
    - System execution mode vs User execution mode
  - Intel defines 4 different privilege levels.

# Privilege levels (II)

- How to scale privileges?
  - Intel offers interrupts
    - Interrupt Driven Operating System
  - When an interrupt/exception happens
    - Hw changes the current privilege level and enables the execution of privilege instructions
  - When the interrupt/exception management ends
    - Hw changes the current privilege level to unable the execution of privilege instructions

# Mechanisms for entering the system

- Exceptions
  - Synchronous, produced by the CPU control unit after terminating the execution of an instruction
- Interrupts
  - Asynchronous, produced by other hardware devices at arbitrary times
- System calls
  - Synchronous: assembly instruction to cause it
    - Trap (in Pentium: INT, sysenter...)
  - Mechanism to request OS services
- All of them are managed through the interrupts vector
  - New architectures implement a fast system call mechanism that skip the interrupts vector: sysenter instruction

# Interrupts Vector

- Pentium
  - IDT: Interrupt Descriptor Table: 256 entries
- Three groups of entries, one for each kind of event:
  - 0 - 31: Exceptions
  - 32 - 47: Masked interrupts
  - 48 - 255: Software interrupts (Traps)

# Initialization

- Each entry in the IDT, identifying an interrupt number, has:
  - A code address
    - Entry point to the routine's code to be executed
  - A privilege level
    - The minimum needed to execute the previous code

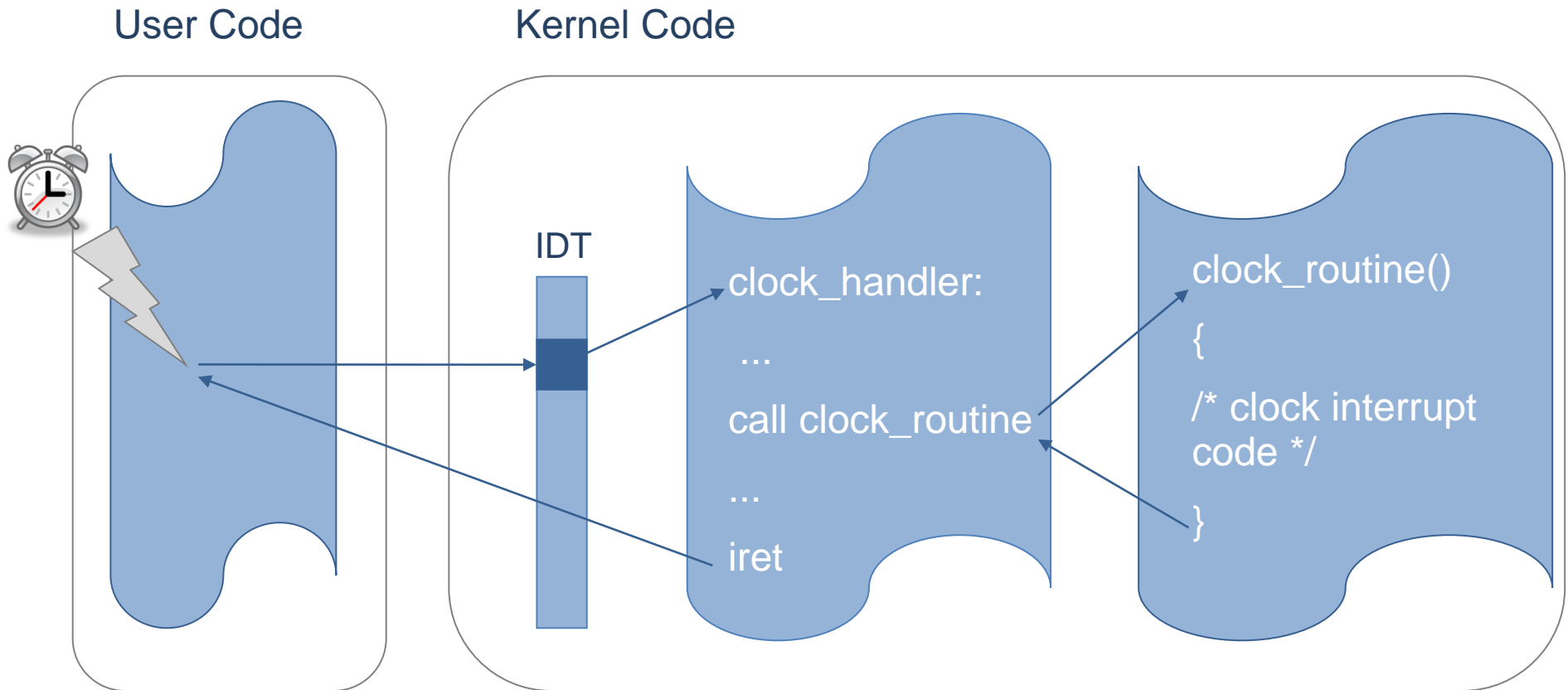


# Management Code


- It could be done in a single routine
  - Divided in two parts: hw context mgmt + solve int.
- Hw context mgmt
  - Entry point handler
  - Basic hardware context management
  - Assembly code
  - Call to a Interrupt Service Routine
- Solve interrupt
  - Interrupt Service Routine
  - High level code (C for example)
  - Specific algorithm for each interrupt

# Example: clock interrupt behavior

Mechanisms for entering the system

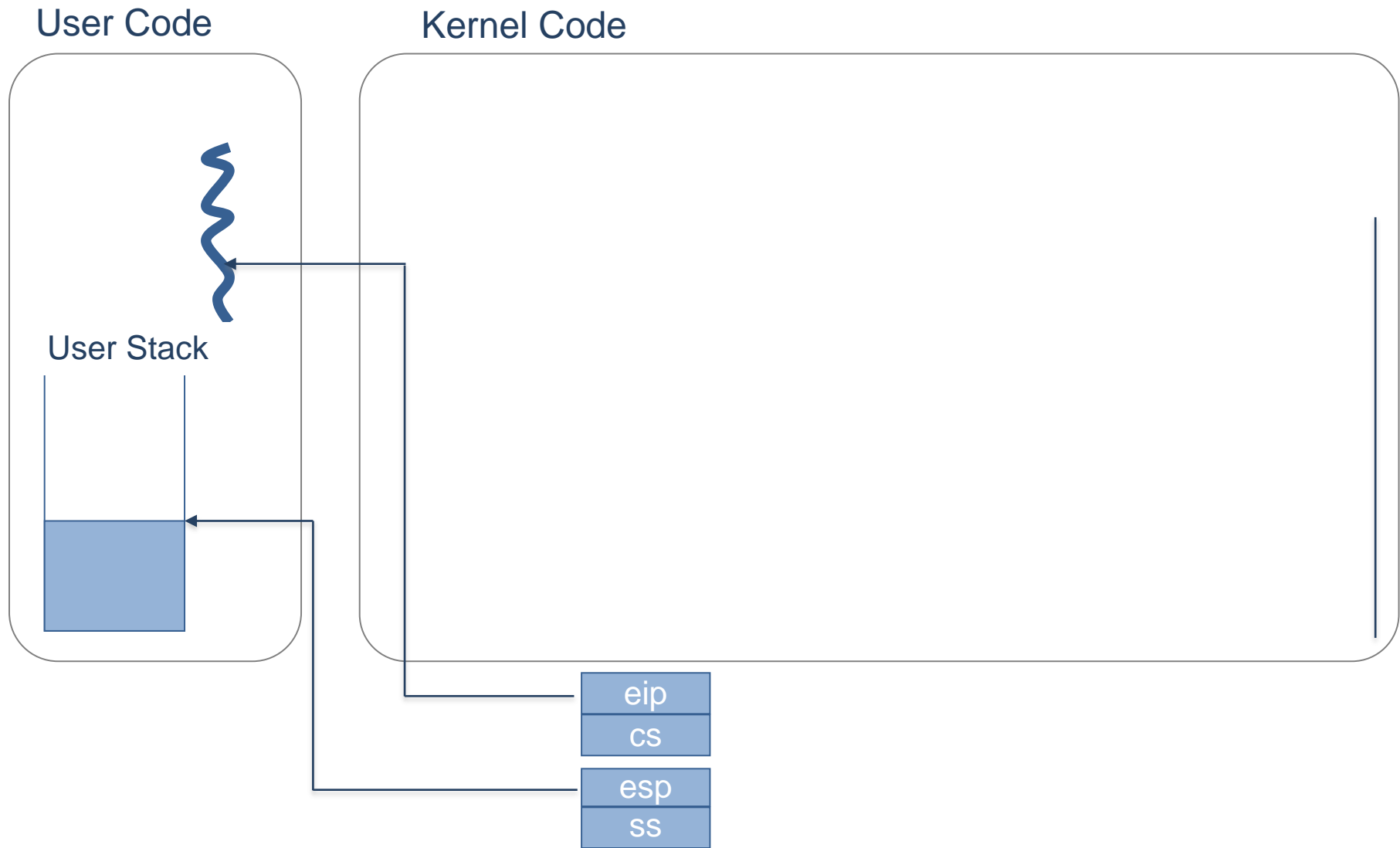


# Procedure for entering the system

- Switch to protected execution mode
    - User Mode → Kernel Mode
  - Save hardware context: CPU registers
    - ss, esp, psw, cs i eip
    - General purpose registers
  - Execute service routine
- 
- The diagram consists of two vertical curly braces on the right side of the list. The top brace spans the first two items of the list (switching to protected execution mode and saving hardware context) and is labeled 'HW'. The bottom brace spans the third item (execute service routine) and is labeled 'handler'.

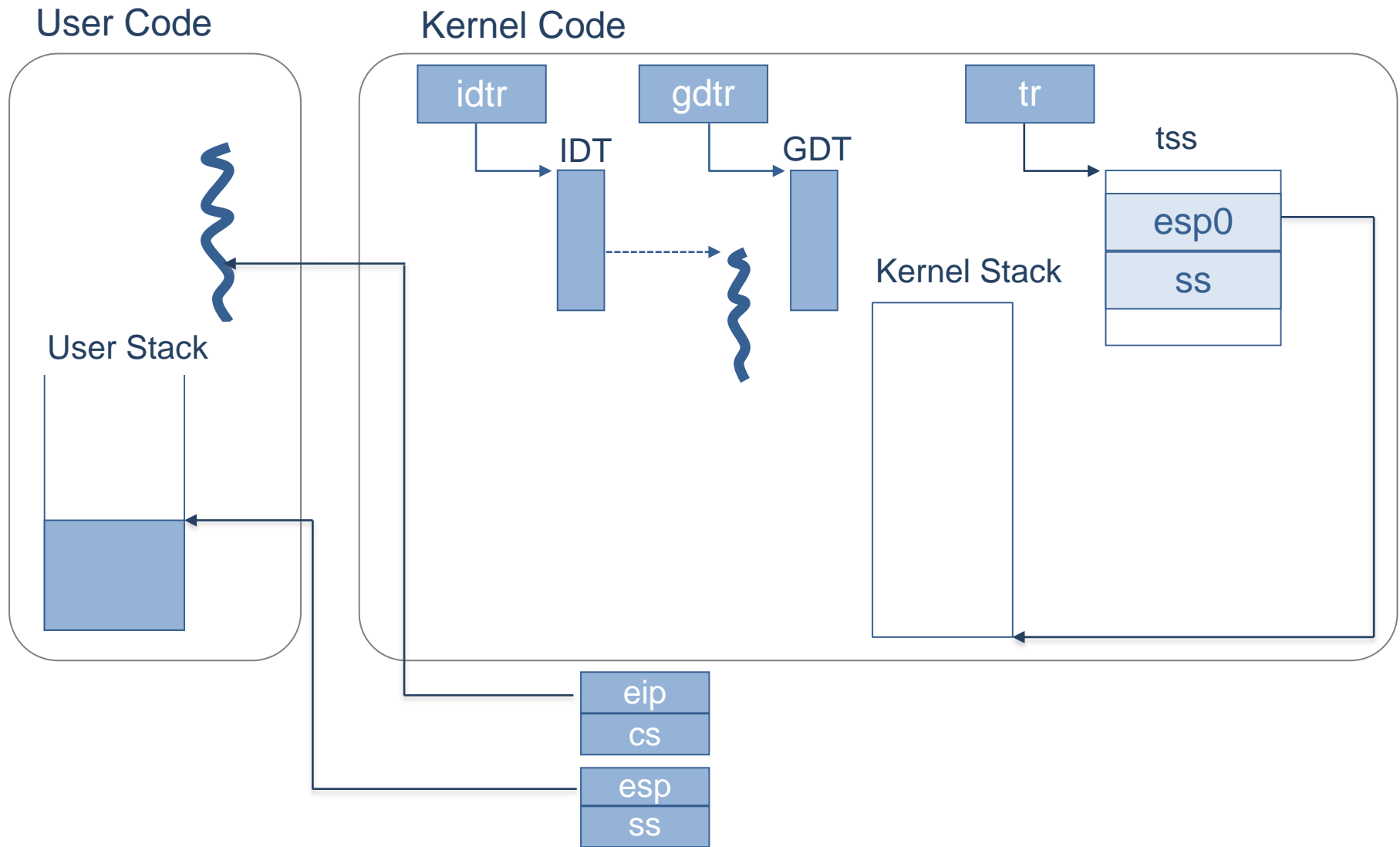
# Procedure for entering the system

Procedure for entering the system



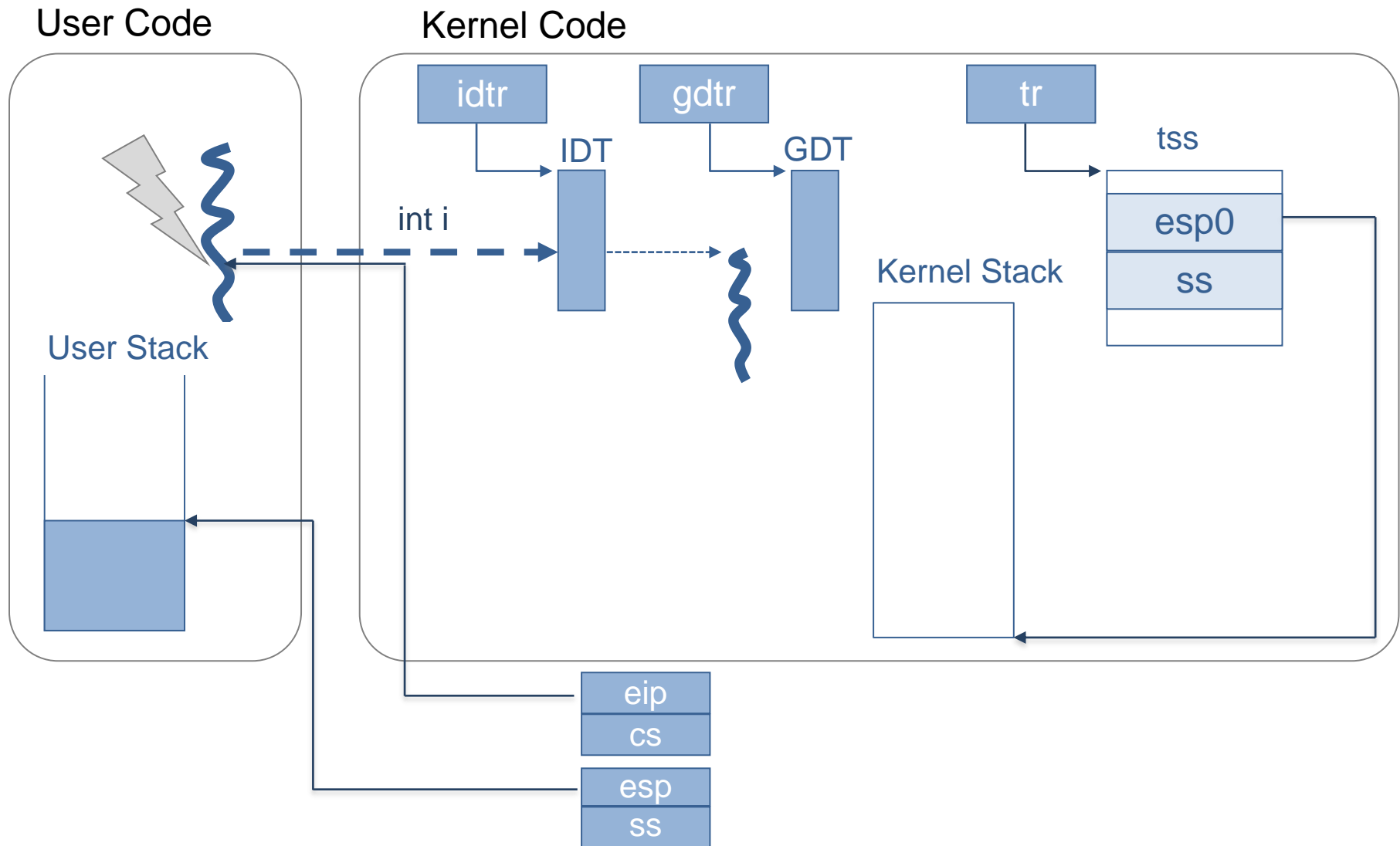
# Procedure for entering the system

Procedure for entering the system



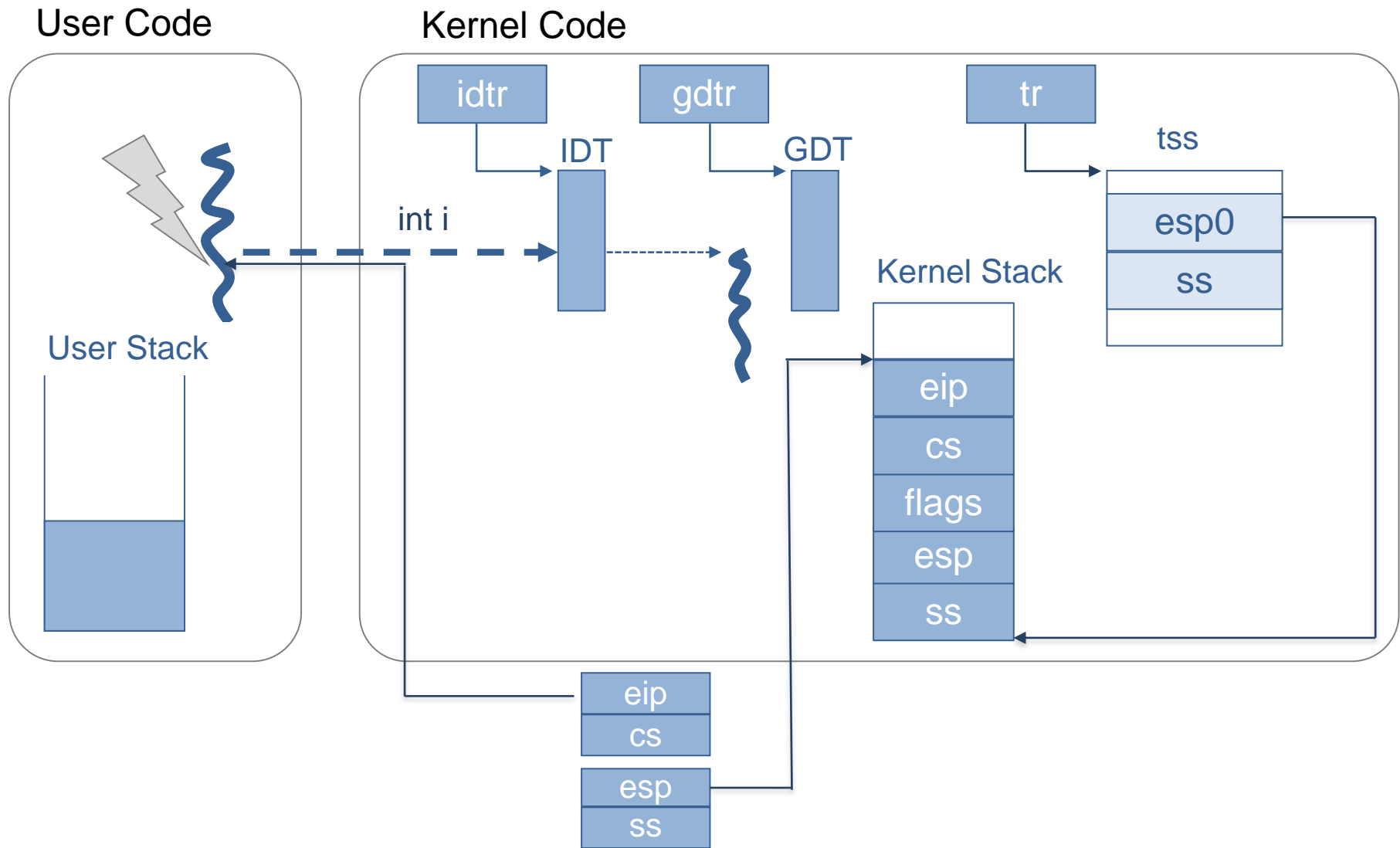
# Procedure for entering the system

Procedure for entering the system



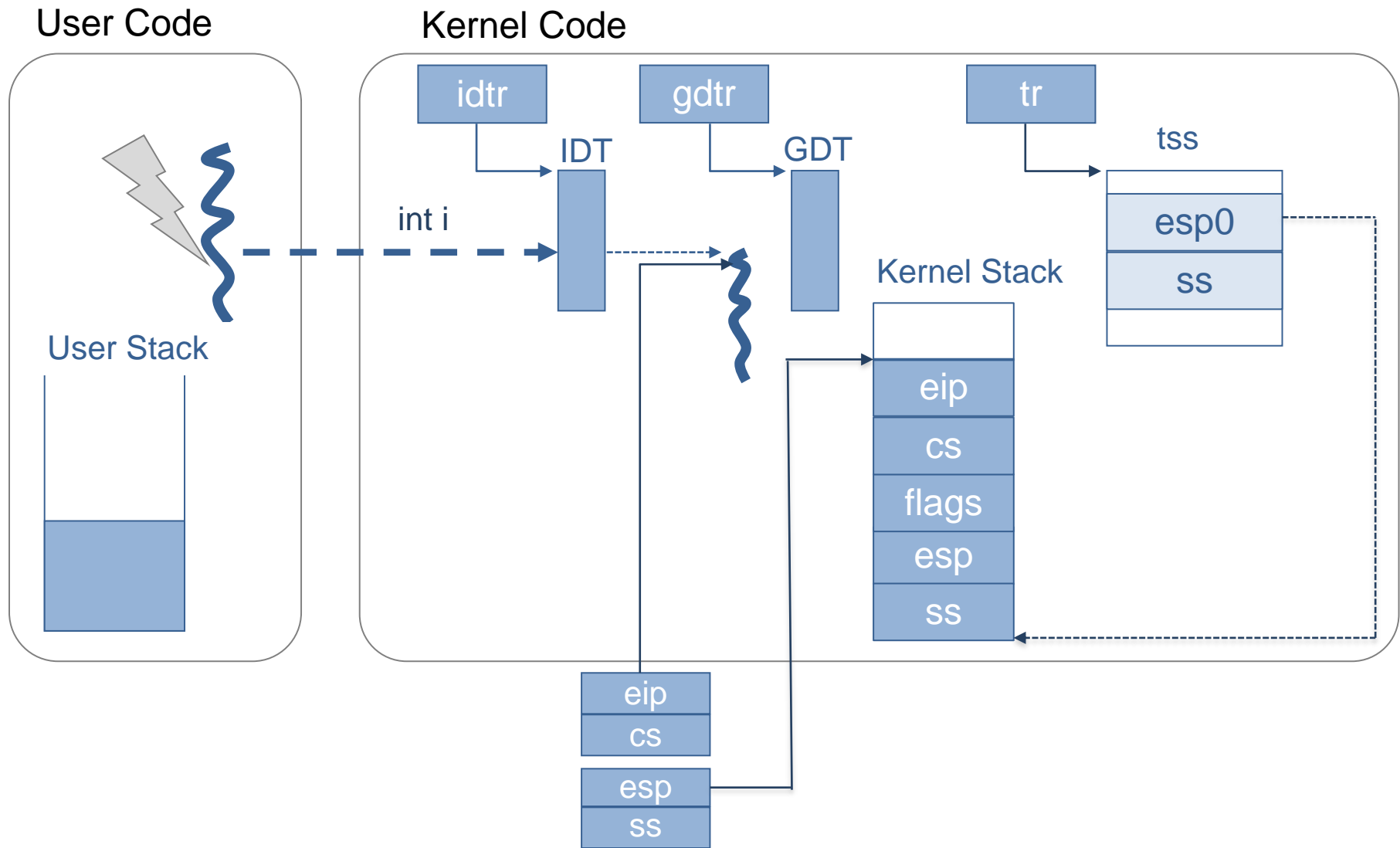
# Procedure for entering the system

Procedure for entering the system



# Procedure for entering the system

Procedure for entering the system





# Procedure to exit the system

- Restore HW context
  - General purpose registers
  - ss, esp, flags, cs, eip
- Switch execution mode
  - Kernel mode → User mode

} handler

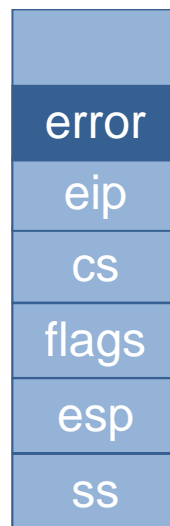
} HW (iret instruction)

Procedure to exit the system

# Exceptions: Stack layout

- There are some exceptions that push a parameter of 4 bytes (a hardware error code) to the kernel stack after entering the system:

Kernel Stack



# Exception's handler

- Save hardware context
- Call exception service routine
- Restore hardware context
- Remove error code (if present) from kernel stack
- Return to user (iret)

# Interrupt's handler

- Similar to exception, but:
  - No hardware error code in kernel stack
  - It is necessary to notify the interrupt controller when the interrupt management finishes
    - Meaning that a new interrupt can be processed
    - End Of Interrupt (EOI)

# Handling system calls

- Why cannot be invoked like a regular user function?
- Which is the mechanism to identify the system call?
- How to pass parameters to the kernel?
- How to get results from the kernel

# System calls: invocation and identification

- Assembly instruction that causes a software generated interrupt
  - `int` assembly instruction (`int idt_entry`)
  - `sysenter` assembly instruction: fast system call mechanism
- An entry point per syscalls?
  - Limitation for the potential number of syscalls
- A single entry point is used for all system calls
  - `int`
    - 0x80 for Linux
    - 0x2e for Windows
  - `sysenter`
    - system call handler @ is kept on a control register: `SYSENTER_EIP_MSR`
- And an extra parameter (EAX) to identify the requested service
- A table is used to translate the user service request to a kernel function to execute

# System calls: parameters and results

- Parameter passing: Stack is NOT shared
  - Linux: syscall handler expects parameters in the registers
    - (first parameter) ebx, ecx, edx, esi, edi, ebp
    - Copy parameters from user stack
  - Windows: Use a register to pass a pointer to parameters
    - EBX
- Returning results:
  - EAX register: contains error code

# System call wrappers

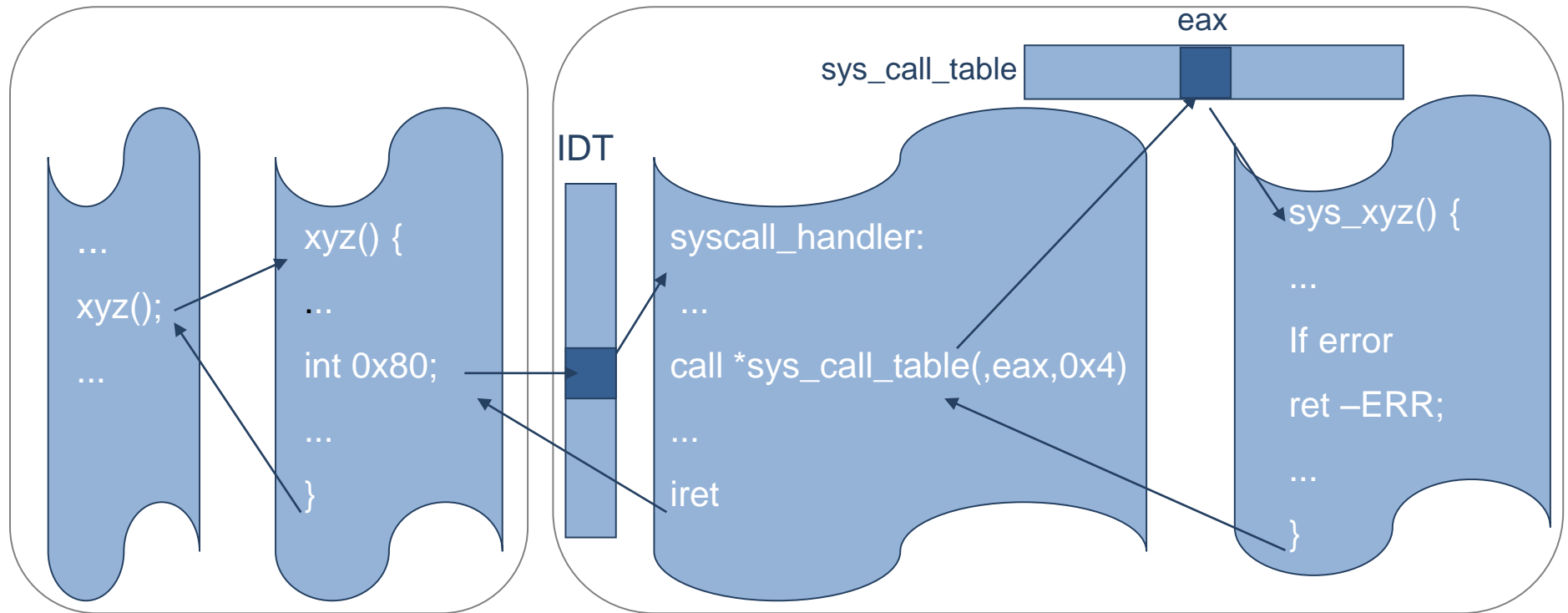
- System must provide the users with an easy and portable way to use them
  - New layer: wrappers
    - wrap all the gory details in a simple function call
- Wrapper responsibilities
  - Invoke the system call handler
    - Responsible for parameter passing
    - Identify the system call requested
    - Generate the trap
  - Return the result to the user code
    - Use errno variable to codify type of error and returns -1 to users



# System call mechanism overview

User Code

Kernel Code



System call invocation in application program

Wrapper for system call

system call handler

system call service routine

System calls

# System call handler

- Save hardware context and prepare parameters for the service routine
  - Linux: stores registers with system call parameters at the top of the kernel stack
  - Windows: copy parameters from the address stored in ebx to the top of the kernel stack
- Execute system call service routine
  - Error checking: system calls identifiers
  - Using `system_call_table`
- Update kernel context with the system call result
- Restore hardware context
- Return to user

# System calls service routines

- Check parameters
  - User code is NOT reliable
    - System MUST validate ALL data provided by users
- Access the process address space (if needed)
- Specific system call code algorithm

# Interrupt Handling Summary

- Save user context
- Restore system context
- Retrieve user parameters [if needed]
- Identify service [if needed]
- Execute service
- Return result [if needed]
- Restore user context