

Operating Systems: Three Easy Steps

Kunwar Shaanjeet Singh Grover

Contents

1	Introduction	3
1.1	Virtualisation	3
1.1.1	CPU Virtualisation	3
1.1.2	Memory Virtualisation	3
1.2	Concurrency	3
1.3	Persistence	3
2	Virtualisation	5
2.1	The Abstraction: The Process	5
2.2	Mechanism: Limited Direct Execution	5
2.2.1	Basic Technique: Limited Direct Execution	6
	Problem #1: Restricted Operations	6
	Problem #2: Switching Between Processes	8

Chapter 1

Introduction

1.1 Virtualisation

The OS takes a **physical** resources (such as the processor , or memory, or a disk) and transforms it into a more general, powerful, and easy-to-use **virtual** for of itself. Thus, we sometimes refer to the operating system as a **virtual machine**. This general technique of transforming is called **virtualisation**.

1.1.1 CPU Virtualisation

Turning a single CPU into a seemingly infinite number of CPUs and thus allowing many programs to seemingly run at once is what we call **virtualizing the CPU**.

1.1.2 Memory Virtualisation

Each process accesses its own **virtual address spaces**, which the OS somehow maps onto the physical memory of the machine. Exactly how this is accomplished is what we study.

1.2 Concurrency

Concurrency is a conceptual term to refer to a host of problems that arise, and must be addressed when working on many things at once (i.e. concurrently) in the same program.

1.3 Persistence

The software in the operating system that usually manages the disk is called the **file system**; it is this responsible for storing any files the user creates in a

reliable and efficient manner on the disks of the system. The file system is the part of OS in charge of managing persistent data. What techniques are needed to do so? What mechanisms and policies are required to do so with a high probability?

Chapter 2

Virtualisation

2.1 The Abstraction: The Process

Process: A running program.

Crux Of The Problem: Although there are only a few physical CPUs available, how can the OS provide the illusion of nearly-endless supply of said CPUs?

The OS creates this illusion by **virtualizing** the CPU. By running one process, then stopping it and running another, and so forth, the OS can promote the illusion that many virtual CPUs exist when in fact there is only one physical CPU (or a few). This technique is called **time-sharing**.

To do this, OS requires some low-level machinery and some high-level intelligence. The low-level machinery is called **mechanisms**. For example, stopping one program and starting another on a given CPU. On top of these mechanisms, there is high-level intelligence call **policies**. Policies are algorithms for making decisions within the OS.

Mechanisms: Provides answer to a *how* question.

Policies: Provides answer to a *which* question.

2.2 Mechanism: Limited Direct Execution

The Crux: How to efficiently virtualize the CPU with control?

The OS must virtualize the CPU in an efficient manner while retaining control over the system. To do so, both hardware and operating-system support will be required. The OS will often use a judicious bit of hardware support in order to accomplish its work effectively.

2.2.1 Basic Technique: Limited Direct Execution

The *direct execution* part of the idea: just run the program directly on the CPU. Thus, when the OS wishes to start a program, it creates a process entry for it in a process list, allocates some memory for it, loads the program code into memory, locates its entry point (`main()`), jumps to it, and starts running the user's code.

This approach gives rise to a few problems:

- How to make sure that the program does not do anything that we don't want it to do while still maintaining efficiency?
- How do we stop a process and switch to another process, i.e. how to implement **time sharing** we require to virtualize the CPU?

Problem #1: Restricted Operations

What if the process wishes to perform some kind of restricted operation, such as issuing an I/O request to a disk, or gaining access to more system resources such as CPU or memory?

The Curx: How to perform restricted operations

A process must be able to perform I/O and some other restricted operations, but without giving the process complete control over the system.

Two different modes:

- **User Mode:** Code that runs in user mode is restricted in what it can do. For example, when running in user mode, a process can't issue I/O requests; doing so would result in the processor raising an exception.
- **Kernel Mode:** Mode the operating system (or kernel) runs in. In this mode, code that runs can do whatever it likes, including privileged operations such as issuing I/O requests and executing restricted instructions.

All modern hardware provides the ability for user programs to perform a **system call**. System calls allow the kernel to carefully expose certain key pieces of functionality to user programs such as accessing the file system, creating and destroying the processes, communicating with other processes, and allocating more memory.

System calls To execute a system call, a program must execute a special **trap** instruction. This instruction jumps into the kernel and raises privilege level to kernel mode; once in the kernel, the system can now perform the restricted operations needed and thus the required work for the calling process. When finished, the OS calls a special **return-from-trap** instruction, which returns into the calling user program while simultaneously reducing the privilege level

back to user mode.

When executing a trap, the caller's registers must be saved in order to return correctly when the OS issues the return-from-trap instruction.

On x86, the processor will push the program counter, flags and a few other registers onto a per-process **kernal stack**; return-from-trap will pop these values from the stack and resume execution.

How does the trap know which code to run inside the OS? Clearly, the calling process can't specify an address to jump to; doing so would allow programs to jump anywhere into the kernel which is not secure. The kernel must control what code executes upon a trap.

The kernel does so by setting up a **trap table** at boot time. When the machine boots up, it does so in kernel mode. One of the first things the OS does is to tell the hardware what code to run when certain exceptional events occur. The OS informs the hardware of these locations of these **trap handlers**, usually with some kind of special instruction. Once the hardware is informed, it remembers the location of these handlers until the machine is rebooted, and thus hardware knows what to do when system calls and other exceptional events take place.

To perform the exact system call, a **system-call number** is usually assigned to each system call. The user code has to place the desired system-call number in a register or at a specific location in the stack; the OS, when handling the system call inside the trap handler, examines the number, ensures it is valid, and, if it is, executes the corresponding code. This level of indirection serves as a form of **protection**; user code cannot specify an exact address to jump to, but rather must request a particular service via a number.

Telling the hardware where the trap tables are is a **privileged** operation, otherwise you could make your own trap tables and compromise the system.

How all this works:

1. At boot time, the kernel initializes the trap table, and the CPU remembers its location for subsequent use. The kernel does so via a privileged instruction.
2. The kernel sets up a few things (allocating memory, etc.) before using a return-from-trap instruction to start the execution of the process. When the process wishes to issue a system call, it traps back into the OS, which handles it and once again returns control via a return-from-trap to the process. The process completes its work and returns from `main()`.

Problem #2: Switching Between Processes

When a process is running on the CPU, this by definition means the OS is *not* running. If OS isn't running how can it change a process? (it can't).

The Crux: How to regain control of the CPU?

How can the operating system **regain control** of the CPU so that it can switch between processes?

Cooperative Approach: Wait for system calls In this approach, the OS regains control of the CPU by waiting for a system call or an illegal operation of some kind to take place, as control is passed to the OS during these exceptions.

Non-Cooperative Approach: The OS takes control If we get stuck in an infinite loop in the cooperative approach, the only way out is to reboot the machine.

The Crux: How to gain control without cooperation

How can the OS gain control of the CPU even if processes are not being cooperative? What can the OS do to ensure a rogue process does not take over the machine?

Timer interrupt: A timer device is programmed to raise an interrupt after a fixed delay; when the interrupt is raised, the currently running process is halted, and a pre-configured **interrupt handler** in the OS runs. Now the OS has regained control of the CPU, and thus can stop the current process and start a different one.

At boot time, the OS starts the interrupt timer. This is a privileged operation.

Saving and Restoring Context The OS has regained control now. It can decide to switch or not. This decision is made by the **scheduler**.

If the decision is to switch, the OS executes a **context switch**. A context switch saves register values of the currently-executing process and restores register values for the soon-to-be-executing process. After this, the OS executes a return-from-trap instruction and executes the new process.

By switching stack pointers (stack pointer is a register), the kernel enters call to the switch code in the context of one process (the interrupted one) and returns in the context of another (the soon-to-be-executing one). By switching stack is actually how the process is switched.

There are two types of register saves/restores that happen here:

- When the timer interrupt occurs. The *user registers* of the running process are implicitly saved by the *hardware*, using the kernel stack of the process.
- When OS performs the context switch. The *kernal registers* are explicitly saved by the *software* (i.e. th OS), but this time into memory in the process structure of the process.