

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	Presistence	3
2	CPU 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
2.3	Scheduling	9
2.3.1	First In, First Out (FIFO)	10
2.3.2	Shortest Job First (SJF)	10
2.3.3	Shortest Time-To-Completion First (STCF)	10
2.3.4	Round Robin	11
	Incorporating I/O	12
2.3.5	The Multi-Level Feedback Queue	12
	Attempt #1: How to Change Priority	13
2.3.6	Problems with current MLFQ	14
2.3.7	Attemp #2: The Priority Boost	14
2.3.8	Attempt #3: Better Accounting	14
3	Memory Virtualization	15
3.1	The Abstraction: Address Space	15
3.2	Mechanisms	15
3.2.1	Dynamic(Hardware-based) Relocation	16
3.2.2	Segmentation	18
	Which Segment Are We Referring To?	18
	Support for Sharing	19
	OS Support	19
3.3	Paging	19

3.3.1	Virtual Address Translation	20
3.3.2	Where Are Page Tables Stored?	21
3.3.3	What is in the Page Table?	21
3.3.4	Paging: Also Too Slow	22
3.3.5	Faster Translations (TLBs)	22
	Who Handles The TLB Miss?	23
	TLB Contents	23
	Context switches with TLB	24
	Replacement Policy	24
3.3.6	Paging: Smaller Tables	25
	Multi-level Page Tables	25
3.4	Swapping	26

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

CPU 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 called **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 **kernel 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 *kernel registers* are explicitly saved by the *software* (i.e. the OS), but this time into memory in the process structure of the process.

More info on the last two points here

2.3 Scheduling

In this section we will understand **policies** that OS scheduling employs.

The Crux: How To Develop Scheduling Policy

How should we develop a basic framework for thinking about scheduling policies? What are the key assumptions? What metrics are important? What basic approaches have been used in earliest of computer systems?

We will make some assumptions and continue to relax them to get to an optimal policy. Assumptions:

1. Each job runs for the same amount of time.
2. All jobs arrive at the same time.
3. Once started, each job runs to completion
4. All jobs only use the CPU (i.e., they perform no I/O)
5. The run-time of each job is known

For now, we will use a single metric to determine optimal policy: **turnaround time**.

$$T_{\text{turnaround}} = T_{\text{completion}} - T_{\text{arrival}}$$

2.3.1 First In, First Out (FIFO)

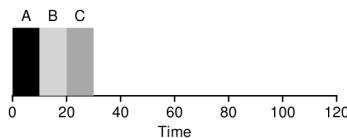


Figure 2.1: FIFO Simple Example

For our current assumptions, this is optimal. But if we relax the assumption that job running time are same, we can get bad turnaround time:

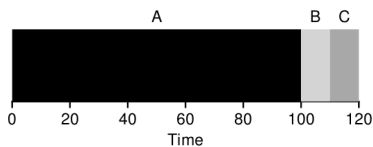


Figure 2.2: Why FIFO is Not that Great

2.3.2 Shortest Job First (SJF)

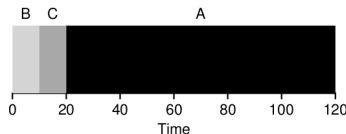


Figure 2.3: SJF Simple Example

It can be proven that under these assumptions, SJF will give the best turnaround time.

Now, let's relax the assumption that all jobs arrive at the same time. Now, we can build a worst case for SJF:

2.3.3 Shortest Time-To-Completion First (STCF)

To address this concern, we need to relax assumption 3 (jobs started, must run to completion). SJF by our definition was a **non-preemptive** scheduler, and thus suffered from the problems in figure 2.4.

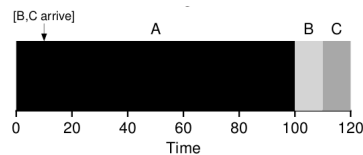


Figure 2.4: SJF With Late Arrivals from B and C

To improve upon that, we add preemption to SJF, known as STCF. Any time a new job enters the system, the STCF scheduler determines which of the remaining jobs has least time left, and schedules that. It can be proven that under the current assumptions, this is the best scheduler (w.r.t turnaround time)

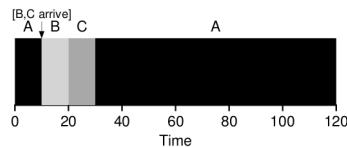


Figure 2.5: STCF Simple Example

New Metric: Response Time: If our only metric was turnaround time and we knew job runtime, STCF would be a great policy. But introduction of time-sharing machines, required interactive performance from system as well. Thus a new metric was needed: **response time**

$$T_{response} = T_{firstrun} - T_{arrival}$$

2.3.4 Round Robin

To improve upon response time, we introduce **Round-Robin (RR)** scheduling. Basic idea: Instead of running a job to completion, RR runs a job for a **time slice** (or **scheduling quantum**) and then switches to the next job in the run queue.

The length of time-slice is obviously a multiple of timer-interrupt.

The length of time slice is critical for RR. The shorter it is, the better performance of RR under response-time metric. However, too small time-slice can result in time taken to context switch to dominate overall performance. Thus, time slices should be made long enough to **amortize** the cost of switching without making it so long that system is no longer responsive.

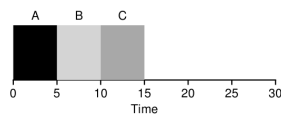


Figure 2.6: SJF (Bad response time)

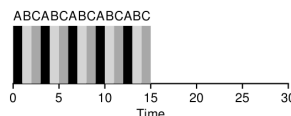


Figure 2.7: RR (Good response time)

RR is one of the *worst* policies for turnaround time because it extends the program to as long as possible.

Incorporating I/O

We now relax assumption that there is no I/O. When a job initiates an I/O request, because the currently-running job won't be using the CPU during the I/O; it is **blocked** waiting for the I/O completion. The CPU is idle during the time of the I/O.

The scheduler incorporates I/O by treating each CPU burst as a job. the scheduler makes sure processes that are "interactive" get run frequently. While those interactive jobs are performing I/O, other CPU-intensive jobs run, thus better utilizing the processor.

2.3.5 The Multi-Level Feedback Queue

In this subsection, we remove the assumption that we know job time and we will describe the most well-known approach to scheduling, **MLFQ**.

The Crus: How To Schedule Without Perfect Knowledge?

How can we design a scheduler that both minimizes response time for interactive jobs while also minimizing turnaround time without a *prior* knowledge of job length?

MLFQ has a number of distinct **queues**, each assigned a different **priority level**. MLFQ uses priorities to decide which job should run at a given time.

Basic rules for MLFQ:

1. If $\text{Priority}(A) > \text{Priority}(B)$, A runs (B doesn't)
2. If $\text{Priority}(A) = \text{Priority}(B)$, A & B run in RR.

Rather than giving a fixed priority to each job, MLFQ *varies* the priority of a job based on its *observed behavior*.

If a job repeatedly relinquishes the CPU while waiting for input from the keyboard, MLFQ will keep its priority high, as this is how an interactive process might behave. If a job uses the CPU intensively for long periods of time, MLFQ will reduce its priority. In this way, MLFQ will try to *learn* about processes as they run, and thus use the *history* of the job to predict its *future* behavior.

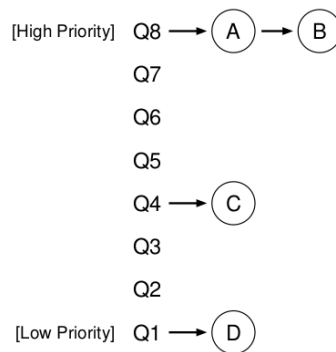


Figure 2.8: MLFQ Example

Attempt #1: How to Change Priority

Keep in mind: Our workload contains a mix of interactive jobs that are short running (and may frequently relinquish the CPU), and some longer running "CPU-bound" jobs that need a lot of CPU time but where the response time isn't important. Here is the first attempt at an algorithm:

3. When a job enters the system, it is placed at highest priority (the topmost queue)
4.
 - If a job uses up an entire time slice while running, its priority is *reduced*.
 - If a job gives up the CPU before the time slice is up, it stays at the *same* priority level.

2.3.6 Problems with current MLFQ

There is a problem of **starvation**: if there are "too many" interactive jobs in the system, they will combine to consume *all* CPU time, and thus long-running jobs will *never* receive any CPU time.

Another problem is that the scheduler can be **gamed**. Gaming a scheduler refers to the idea to tricking the scheduler in giving more time than the job's share. Before the time-slice is over, the job could issue an I/O and remain at the same priority. Doing so, allows the program to gain a higher percentage of CPU time.

2.3.7 Attempt #2: The Priority Boost

We add a new rule to handle the flaws.

5. After some time period S , move all the jobs in the system to the topmost queue.

This solves the problem of starvation. S has to be set properly.

2.3.8 Attempt #3: Better Accounting

We rewrite Rule 4 to make it anti-gaming.

4. Once a job uses its time allotment at a given level (regardless of how many times it has given up the CPU), its priority is reduced.

MLFQ: Summary

1. If $\text{Priority}(A) > \text{Priority}(B)$, A runs (B doesn't)
2. If $\text{Priority}(A) = \text{Priority}(B)$, A & B run in round-robin fashion using the time slice of the given queue.
3. When a job enters the system, it is placed at the highest priority.
4. Once a job uses up its time allotment at a given level, its priority is reduced.
5. After some period S , move all jobs to the topmost queue.

Chapter 3

Memory Virtualization

3.1 The Abstraction: Address Space

Address Space: The running program's view of memory in the system

The Crux: How to Virtualize Memory?

How can the OS build this abstraction of a private, potentially large address space for multiple running processes (all sharing memory) on top of a single, physical memory?

Goals of OS while virtualizing memory:

1. Transparency: How OS implements virtual memory should be invisible to the running program.
2. Efficiency
3. Protection: A process should not be able to access memory out of its virtual memory

3.2 Mechanisms

The Crux: How to Efficiently And Flexibly Virtualize Memory

How can we build an efficient virtualization of memory? How do we maintain control over the locations application can access? how do we do all of this efficiently?

How it will work:

- Hardware provides mechanisms for **hardware-based address translation**. Hardware only provides the low-level mechanisms.
- OS will use these mechanisms to manage memory, keeping track of free memory and maintain control.

Assumptions we make for now:

- Address space must be placed *contiguously* in physical memory.
- Size of address space is less than size of physical memory.

We will continue to relax these assumptions as we develop a better model. The current model will be laughable at best.

What the process should see:

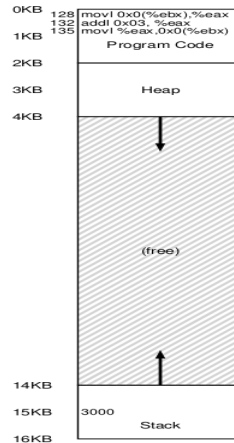


Figure 3.1: A process and its address space

What should actually go in memory:

The program thinks its address space starts at address 0, while it actually starts at 32KB in the physical memory.

3.2.1 Dynamic(Hardware-based) Relocation

We discuss the first incarnation of address translation: **base and bounds**.

Specifically, we'll need two hardware registers within each CPU: one called the **base** register, and the other the **bounds** register. Using base and bounds, we can place the address space anywhere we like in the physical memory and ensure that it only has access to its own address space.

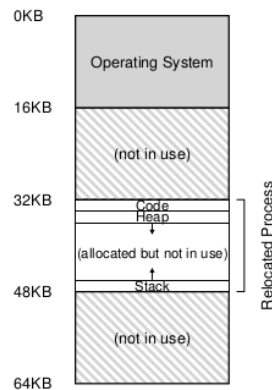


Figure 3.2: Physical Memory with the process

In this setup, each program is written and compiled as if it is loaded at address zero. However, when running the OS decides where in physical memory to place the program and sets the base register to that value. Now, the physical address can be calculated as:

$$\text{physical address} = \text{virtual address} + \text{base}$$

Memory references generated by the process are **virtual addresses**. The hardware turns these virtual addresses into **physical addresses**. This technique of transforming a virtual address is exactly what is referred to as **address translation**. The bounds register is used to check if the virtual address lies in the virtual memory or not.

Since this relocation happens at runtime and we can move address spaces after process has started running, this technique is referred to as **dynamic relocation**.

Note that base and bounds registers are hardware structures kept on the CPU chip. This part of processor which helps with address translation is called **Memory Management Unit (MMU)**. We will continue to add more circuitry to MMU.

Hardware Support Summary till now

- **Privileged mode:** Needed to prevent user-mode processes from executing privileged instructions
- **Base/bounds registers:** pair of registers per CPU to support address translation and bound checks

- **Privileged instructions to update base/bounds**
- **Privileged instructions to register exception handlers:** How to handle exceptions must be told by the OS
- **Ability to raise exceptions**

OS Requirement Summary for memory

- **Memory Management:** Need to allocate memory for new processes, reclaim memory from terminated processes, manage memory via **free list**.
- **Base/bounds management:** Must set base/bounds properly upon context switch.
- **Exception handling:** Code to run when exceptions arise.

Base and bounds virtualization is quite *efficient* but it has a **lot** of problems. As can be seen in Figure 3.2 all of the space between the stack and heap is wasted. This is called **internal fragmentation**.

3.2.2 Segmentation

The Crux: How To Support A Large Address Space

How do we support a large address space with (potentially) a lot of free space between the stack and heap?

Basic idea → **Segmentation**. Instead of having just one base and bounds pair in our MMU, why not have a base and bounds pair per logical segment of the address space? In our address space, we have three logically-different segments: code, stack, and heap. In segmentation, we place each one of those segments in different parts of memory.

Which Segment Are We Referring To?

Two approaches:

1. **Explicit** approach: Chop up the address space into segments based on few bits of the virtual address.

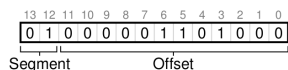


Figure 3.3: Chopping virtual address

2. **Implicitly** approach: The hardware determines the segment by noticing how the address was formed.

Since our memory can grow forward or backward, we add another piece of information in hardware that whether the segment grows positively or negatively.

Support for Sharing

To save memory, sometimes it is useful to **share** certain memory segments between address spaces. In particular, **code sharing** is common nowadays.

To support this, we need to add more support for hardware in form of **protection bits**. Basic support adds a few bits per segment, indicating read/write/execute permissions. Trying to do something out of permission should raise an exception.

OS Support

1. The OS has to save all base and bound register when a context switch occurs.
2. OS interaction when segments grow: if `malloc()` is called, in some cases the heap may be able to service the object. If it can, find free space for the object and return a pointer to it. If the heap needs to grow, call a system call to grow the heap.
3. Managing free space: When a new address is created, the OS has to be able to find space in physical memory for its segments. The general problem that arises is that physical memory soon becomes full of holes and to allocate a contiguous block of memory, you will have to move some blocks around. This problem is called **external fragmentation**.

The main problem right now is using a variable sized chunks.

3.3 Paging

Chop up space into *fixed-size* pieces → **paging**

The Curx: How To Virtualize Memory With Pages

How can we virtualize memory with pages, so as to avoid the problems faced by segmentation? What are the techniques and how do we make them work efficiently?

Paging, has a number of advantages over our previous approaches:

- Probably the most important improvement will be *flexibility*: with a fully-developed paging approach, the system will be able to support the abstraction of an address space effectively, regardless of how a process uses the address space. For example, no assumptions in the way stack and heap grow.
- Another advantage is the *simplicity* of free-space management that paging affords. For example, when the OS wishes to allocate lets say 64-byte address space into the physical memory, it just finds four free pages for this. OS keeps some kind of **free list** of all free pages for this.

To record where each virtual page of the address space is placed in physical memory, the operating system usually keeps a *per-process* data structure known as **page table**. The major role of the page table is to store **address translations** for each of the virtual pages of the address space, thus letting us know where in physical memory each page resides.

The page table is a *per-process* data structure (most page table structures we do here are like this, an exception is **inverted page table**). If another process wanted to run, we would have to change our page table.

3.3.1 Virtual Address Translation

Let's imagine the process with a tiny address space (64 bytes) is performing a memory access:

```
movl <virtual address>, \%eax
```

To **translate** this virtual address that the process generated, we have to first split it into two components: the **virtual page number (VPN)**, and the **offset** within this page. For this example, we need need 6 bits total as $2^6 = 64$. Because we know the page size (16 bytes), we divide the virtual address as follows:

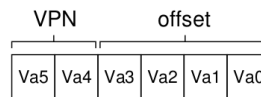


Figure 3.4: Virtual Address of the example

Lets say that we have our virtual address as 21, the VPN is 01, lets say it maps to the 7th **physical frame number (PFN)** also called the **physical page number (PPM)**. Thus we can translate this address by replacing the VPN with the PFN and then issue the load to physical memory.

Note that the offset stays the same, because the offset just tells us which byte *within* the page we want.

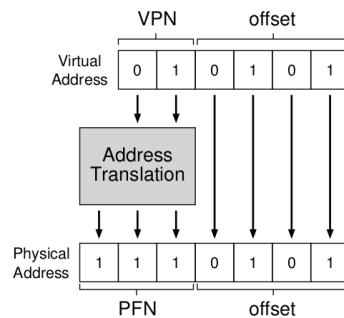


Figure 3.5: The Address Translation Process

3.3.2 Where Are Page Tables Stored?

Page tables can get terribly large, much bigger than the small segment table or base/bounds pair we discussed previously. For example, for a 32-bit address space, with 4KB pages, we have a 20-bit VPN and 12-bit offset. A 20-bit VPN implies that there are 2^{20} translations that the OS would have to manage for each process. Assuming we need 4 bytes per **page table entry (PTE)** to hold the physical translation plus any other useful info, we get 4MB of memory needed for each page table, i.e. we need 4MB for each process. If we have a 100 processes running at a time, the Page Tables are using 400MB, which is too much. This will be even worse for 64-bit address space. Since they are so big, we don't keep any special onchip-hardware in MMU to store. We store it in the memory.

3.3.3 What is in the Page Table?

The page table is just a data structure used to map virtual addresses to physical addresses. For now, we use a **linear page table**, which is just an array. Later we'll use more complex structures to reduce memory.

We add the following things to a PTE:

1. A **valid bit** to indicate if the particular translation is valid. For example, on a running code we have stack and heap on opposite ends. The space inbetween them will be invalid. Accessing invalid positions will generate a trap to OS which will say GOODBYE OFFENDING PROCESS.
2. **protection bits** indicating read/write/execute permissions.
3. A **present bit** indicates whether this page is **swapped out** (Will study this further)
4. A **dirty bit** indicating if the page has been modified since the last time it was brought into memory.

5. **reference bit** Useful in determining which pages are popular. Such knowledge is critical during **page replacement** (Further)

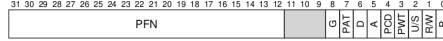


Figure 3.6: An x86 Page Table Entry (PTE)

3.3.4 Paging: Also Too Slow

The OS first goes to the memory and looks at the page table, brings back the translation, then again goes to memory to access the physical address mapped to the virtual address. Here is code for how it is done:

```

1 // Extract the VPN from the virtual address
2 VPN = (VirtualAddress & VPN_MASK) >> SHIFT
3
4 // Form the address of the page-table entry (PTE)
5 PTEAddr = PTBR + (VPN * sizeof(PTE))
6
7 // Fetch the PTE
8 PTE = AccessMemory(PTEAddr)
9
10 // Check if process can access the page
11 if (PTE.Valid == False)
12     RaiseException(SEGMENTATION_FAULT)
13 else if (CanAccess(PTE.ProtectBits) == False)
14     RaiseException(PROTECTION_FAULT)
15 else
16     // Access is OK: form physical address and fetch it
17     offset = VirtualAddress & OFFSET_MASK
18     PhysAddr = (PTE.PFN << PFN_SHIFT) | offset
19     Register = AccessMemory(PhysAddr)

```

Figure 3.7: Accessing Memory with Paging

We have extra memory references which likely slow down the process by a factor of two or more.

Problem with paging → if not designed carefully, it will cause the system to run too slow as well as take too much memory. We need to overcome these two problems now.

3.3.5 Faster Translations (TLBs)

The Crux: How To Speed Up Address Translation

How can we speed up address translation, and generally avoid the extra memory reference that paging seems to require? What hardware support is required? What OS involvement is needed?

To speed up address translation, we are going to add **translation- lookaside buffer or TLB**. A TLB is part of the chip's **memory- management unit (MMU)**. It is basically **cache** for popular virtual-to-physical translations.

The algorithm hardware follows is like this:

- Extract the VPN from virtual address
- If the VPN is present in the TLB then it is a **TLB hit**. Extract the PFN and concatenate it with the offset.
- If it is a **TLB miss**, hardware accesses the page table to find the translation, assuming that the virtual reference generated by the process is valid and accessible, updates the TLB with the translation. A TLB miss is expensive because of the extra memory reference.
- Finally, once the translation is found in the TLB, the memory reference is processed quickly.

The TLB like all caches is based on spatial and temporal locality.

Who Handles The TLB Miss?

TLB miss can be handled by the hardware or the OS.

In the olden days, the hardware had **CISC** architecture and handled the TLB miss entirely. To do this, the hardware has to know exactly *where* the page tables are located in memory as well as their *exact format*. On a miss, the hardware would find the page-table entry and extract the desired translation, update the TLB with the translation and retry the instruction. Intel x86 architecture uses **hardware-managed page TLBs** with a fixed **multi-level page table** (will study further)

More modern architectures with **RISC** style have what is known as a **software-managed TLB**. On a TLB miss, the hardware simply raises an exception, which pauses the current instruction stream, raises the privilege level to kernel mode and jumps to a **trap handler**. The trap handler handles the miss. Main reason of using OS to handle miss is more *flexibility*.

TLB Contents

A TLB has working similar to cache. But what is the data in it? A TLB entry might look like this:

VPN | PFN | other bits

The interesting part is the "other bits". For example, the TLB commonly has a **valid** bit, which says whether the entry has a valid translation or not. Also

common are **protection** bits, which determine the permissions the process has for the memory. There might be other fields, including an **address-space identifier**, a **dirty bit**, and so forth.

Context switches with TLB

The Crux: How To Manage TLB Contents On A Context Switch

When context-switching between processes, the translations in the TLB for the last process are not meaningful to the about-to-be-run process. What should the hardware or OS do in order to solve this problem?

One way would be to flush the TLB on context switches, thus emptying it before running the new process. However this is costly if the processes are switched frequently.

To reduce this overhead, some systems add hardware support to enable sharing of the TLB across context switches. In particular some systems provide an **address space identifier (ASID)** field in the TLB which can be thought of as a **process identifier (PID)**, but usually takes fewer bits. Here is a depiction of TLB with the added ASID field:

VPN	PFN	valid	prot	ASID
10	100	1	rwX	1
—	—	0	—	—
10	170	1	rwX	2
—	—	0	—	—

Figure 3.8: TLD entry with ASID

Now the TLB can hold translations from different processes at the same time without any confusion. The hardware also needs to know which process is currently running in order to perform translations, and thus the OS must, on a context switch, set some privileged register to the ASID of the current process.

Replacement Policy

The replacement policies are same as **cache replacement**. Duh, its a cache. What were you expecting? Anyway, this part is covered in swapping.

3.3.6 Paging: Smaller Tables

The Crux: How To Make Page Tables Smaller

Simple array-based page tables are too big. How can we make page tables smaller? What are the key ideas? What inefficiencies arise as a result of these new data structures?

Some solutions with their problems:

- **Bigger Pages:** Can lead to internal fragmentation due to pages not being used completely.
- **Paging + Segments:** Divide segments into pages. Doesn't really solve the main problem of segmentation i.e. external fragmentation.

Multi-level Page Tables

Get rid of all the invalid regions in the page table instead of keeping them all in memory.

We turn the linear page table into something like a tree. The basic idea is simple, chop up the page table into page-sized units; then if an entire page of page-table entries (PTEs) is invalid, don't allocate that page of the page table at all. To track whether a page of the page table is valid, use a new structure called the **page directory**. the page tells you where a page of the page table is or that the entire page of the page table contains no valid pages.

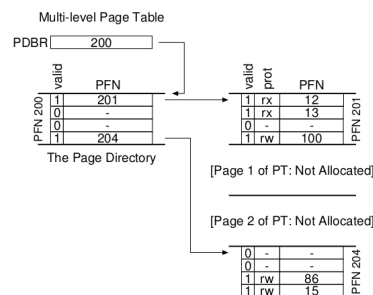


Figure 3.9: 2-Level Page Table

Advantages of Multi-level page tables:

- They only allocate page-table space in proportion to the amount of address space you are using; thus it is good for sparse address spaces

- If carefully constructed, each portion fits neatly within a page, making it easier to manage memory.

It should be noted that there is a cost to multi-level tables; for the 2-level page table, we required two loads from memory to get the right translation information from the page table (one for the page directory and one for the PTE itself). Smaller tables give higher TLB miss cost but reduce space. This is a space-time tradeoff. The cost in a TLB hit is the same.

We don't have to use only 2-level page tables. For larger physical memory, we need to use higher level page tables to get smaller tables. We create a page directory of page directories and so on. Levels can be decided based on how much reduction we need.

New code for Address translation:

```

1  VPN = (VirtualAddress & VPN_MASK) >> SHIFT
2  (Success, TlbEntry) = TLB_Lookup(VPN)
3  if (Success == True) // TLB Hit
4      if (CanAccess(TlbEntry.ProtectBits) == True)
5          Offset = VirtualAddress & OFFSET_MASK
6          PhysAddr = (TlbEntry.PFN << SHIFT) | Offset
7          Register = AccessMemory(PhysAddr)
8      else
9          RaiseException(PROTECTION_FAULT)
10 else // TLB Miss
11     // first, get page directory entry
12     PDIndex = (VPN & PD_MASK) >> PD_SHIFT
13     PDEAddr = PDBR + (PDIndex * sizeof(PDE))
14     PDE = AccessMemory(PDEAddr)
15     if (PDE.Valid == False)
16         RaiseException(SEGMENTATION_FAULT)
17     else
18         // PDE is valid: now fetch PTE from page table
19         PTIndex = (VPN & PT_MASK) >> PT_SHIFT
20         PTEAddr = (PDE.PFN << SHIFT) + (PTIndex * sizeof(PTE))
21         PTE = AccessMemory(PTEAddr)
22         if (PTE.Valid == False)
23             RaiseException(SEGMENTATION_FAULT)
24         else if (CanAccess(PTE.ProtectBits) == False)
25             RaiseException(PROTECTION_FAULT)
26         else
27             TLB_Insert(VPN, PTE.PFN, PTE.ProtectBits)
28             RetryInstruction()

```

Figure 3.10: Multi-level Page Table Control Flow

As it can be seen, only the cost of TLB miss has been increased and not that of a TLB hit.

3.4 Swapping

Now we relax the assumption that our address space is unrealistically small and fits into physical memory,

The Crux: How To Go Beyond Physical Memory

How can the OS make use of a larger, slower device to transparently provide the illusion of a large virtual address space?

First, we need to reserve some space for swapping. we generally refer to such space as **swap space**. Thus, we simply assume that the OS can read from and write to the swap space, in page-sized units. OS needs to remember the **disk address** of a given page to accomplish this.

Now that we have some space on disk for swapping, we need to add machinery in the system to support swapping pages to and from the disk. We assume that we have a system with a hardware-managed TLB (like x86 has).

We need to add more machinery to the **page table entry (PTE)**. It must now have a **preset bit**, which indicates whether the page is in memory. In the case that the hardware looks in the TLB and the page is *not present*