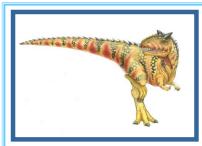


Chapter 9: Main Memory



Chapter 9: Memory Management Strategies

- Background
- Contiguous Memory Allocation
- Segmentation
- Paging
- Structure of the Page Table
- Examples: the Intel 32 and 64-bit Architectures and ARM Architecture



Objectives

- Explain the difference between a logical and a physical address and the role of the memory management unit (MMU) in address translation.
- Dynamic storage-allocation problem - apply **first-fit**, **best-fit**, and **worst-fit** strategies for allocating memory contiguously.
- Explain the distinction between **internal** and **external fragmentation**.
- Translate logical to physical addresses in a paging system that includes a translation look-aside buffer (TLB).
- Describe hierarchical paging, and hashed paging
- Describe address translation for IA-32, x86-64, and ARMv8 architectures



Background

- Main memory is central to the operation of computer systems
- Memory consists of a large array of **bytes**, each with its own address – **byte-addressable**
- Program must be brought (from disk or other secondary storage) into memory and placed within a process for it to run, i.e., PCB points to the address space of the process
- A typical instruction-execution cycle include
 - CPU first fetches an instruction from memory (or cache)
 - The instruction is then decoded and may cause operands to be fetched from memory.
 - After the instruction executed on the operands, results may be stored back in memory.
- The **memory management unit** or **MMU** only sees a stream of **addresses**
 - It does not know how they are generated, could be by the instruction counter, indexing, indirection, literal addresses, and so on or what they are for (instructions or data).
 - Accordingly, MMU is interested only in the sequence of memory addresses generated by a running program.



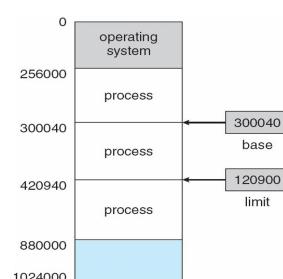
Background (Cont.)

- Main memory (including cache) and registers are the only storage CPU can access directly; in another word, CPU cannot access secondary storage such as disk directly, but through I/O controllers
- Memory management unit or MMU only sees a stream of **addresses + read requests**, or **address + data and write requests**
 - Registers can be accessed in one CPU clock (cycle)
 - Accessing main memory (through memory bus) may take many CPU clocks, causing a **memory stall**, when it does not yet have the data required to complete the instruction it is executing
 - Cache sits between main memory and CPU registers, residing on CPU chips for fast access
- Memory protection is required to ensure correct operation
 - We must protect the operating system memory space from being accessed by user processes, as well as protect user processes from one another
 - This protection must be provided by the hardware for performance – performance or speed



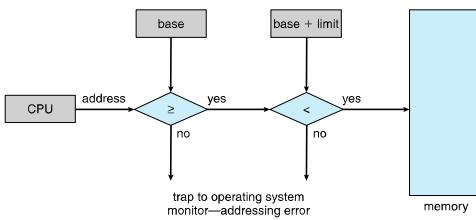
Per-process memory space separated from each other

- Separate per-process memory space protecting the processes from each other is fundamental to have multiple processes loaded in memory for concurrent execution – multiprogramming systems
- A pair of **base** and **limit registers** define the legal range of a process address space
 - For example, if the base register holds 300040 and the limit register is 120900, then the program can legally access all addresses from 300040 through 420939 (inclusive)



Hardware Address Protection

- Protection of memory space is accomplished by having CPU compare every address generated in **user mode** with these two registers
- The **base** and **limit** registers can be loaded only by the operating system, which uses a special privileged instruction (executed only in the **kernel mode**)
- The operating system, executing in **kernel mode**, however, is given unrestricted access to both operating-system memory and users' memory



Operating System Concepts – 10th Edition

9.7

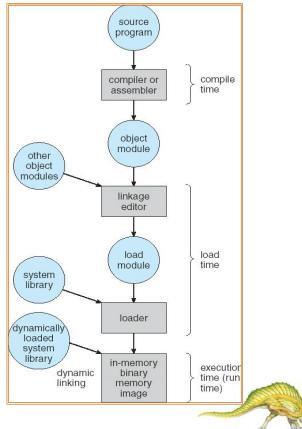
Operating System Concepts – 10th Edition

9.8

Multipstep Processing of a User Program

Address binding of instructions and data to memory addresses can happen in three different stages

- **Compile time:** If memory location is known a priori, **absolute code** can be generated; must recompile code if starting location changes (e.g., “gcc”). MS-DOS uses this
- **Link or Load time:** Compiler must generate **relocatable code** if memory location is not known at compile time (e.g. Unix “ld” does link). The binding is delayed until load time. If the starting address changes, we need only **reload the user code** to incorporate this changed value
- **Execution time:** Binding delayed until run time if the process can be moved during its execution from one memory segment to another. This needs **hardware** and **operating system support** for address maps (e.g., **base** and **limit registers**), e.g., dynamic libs. **Most general-purpose operating systems use this method.**



Operating System Concepts – 10th Edition

9.9

Operating System Concepts – 10th Edition

9.10

Address Binding 地址绑定

- Usually, a program resides on a disk as a **binary executable file**, which must be brought into memory and placed within a process (part of its address space) for it to run
- Most systems allow a user process to reside in any part of the physical memory. Although the computer may start at 00000, the first address of a user process need not be 00000
- A user program goes through several steps - some of which may be optional before being executed. Addresses may be represented in different ways during these steps:
 - Source code addresses usually symbolic – the variable count
 - A compiler typically **binds** these symbolic addresses to relocatable addresses, such as “14 bytes from beginning of this module”
 - **Linker or loader** in turn **binds the relocatable addresses to absolute addresses**, i.e., 74014
 - Each binding is essentially a mapping from one address space to another



Operating System Concepts – 10th Edition

9.8

逻辑地址空间 vs. 物理地址空间

- Recall: Address Space
 - All the addresses a process can “touch”
 - Each process has its own unique address space, different from kernel address space
- Consequently, there are two views of memory
 - **Logical address** – generated by the CPU; also referred to as **virtual address**
 - **Physical address** – address seen by the memory
 - The translation is done by **memory management unit or MMU** – hardware device
- Translation makes it much easier to implement protection
 - To ensure a process can not access another process's address space,
- **Logical and physical addresses are the same if compile-time and load-time address-binding schemes are used (in old days); logical (virtual) and physical addresses differ only in execution-time address-binding scheme**
 - **Logical address space** is the set of all logical addresses generated by a program
 - **Physical address space** is the set of all physical addresses generated by a program

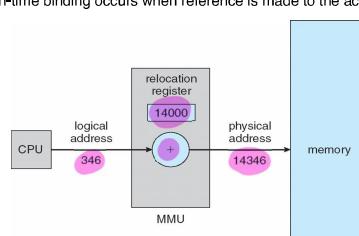


Operating System Concepts – 10th Edition

9.11

Memory-Management Unit (MMU)

- MMU is a **hardware mechanism** that at runtime maps virtual address to physical address
- There are many different mapping methods, covered in the rest of this chapter
- To start, consider a simple scheme where the value in the **relocation register** is added to every address generated by a user process at the time it is sent to memory
 - **User programs never accesses the real physical addresses**
- The user program deals with **logical addresses**; it never sees the **real physical addresses**
 - Execution-time binding occurs when reference is made to the actual location in memory



Operating System Concepts – 10th Edition

9.12





Contiguous Allocation 连续分配

- Contiguous allocation is one of the early memory allocation methods
- Main memory is usually divided into two partitions: one for operating system and one for user processes
- The operating system can be placed in either low memory addresses or high memory addresses, depending on many factors, such as the location of the interrupt vector
- Many operating systems (including Linux and Windows) place the operating system in high memory
- In a multi-programmed operating system, several user processes reside in memory at the same time. In contiguous memory allocation, each process contained in a single section of memory that is contiguous to the section containing the next process.

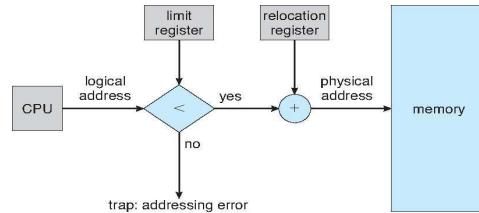
Operating System Concepts – 10th Edition

9.13



Hardware Support for Relocation and Limit Registers

- Relocation register and limit register are used to protect user processes from each other, and from changing operating-system code and data
- Relocation register contains value of the smallest physical address
- Limit register contains range of logical addresses – each logical address must be less than the limit register
- MMU maps logical address dynamically

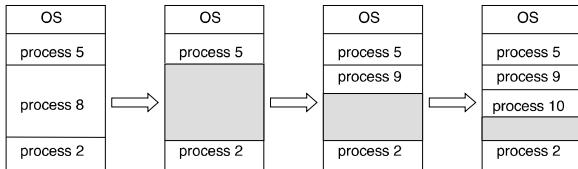
Operating System Concepts – 10th Edition

9.14



Contiguous Allocation (Cont.)

- Multiple-partition allocations 多分区
 - Degree of multiprogramming is bounded by number of partitions
 - Variable-partition sizes (sized to a given process' needs)
 - Hole** – block of available memory; holes of various size are scattered throughout memory
 - When a process arrives, it is allocated memory from a hole large enough to accommodate it
 - Process exiting frees its partition, adjacent free partitions combined
 - Operating system maintains information about a) allocated partitions; b) free partitions (hole)

Operating System Concepts – 10th Edition

9.15



Dynamic Storage-Allocation Problem

How to satisfy a request of (variable) size n from a list of free holes?

- First-fit:** Allocate the *first* hole that is big enough
- Best-fit:** Allocate the *smallest* hole that is big enough
 - Must search entire list, unless ordered by size
 - Produces the smallest leftover hole – intended not use it
- Worst-fit:** Allocate the *largest* hole
 - Must also search entire list
 - Produces the largest leftover hole – intended to reuse the remaining hole

Experiments have shown that both first fit and best fit are better than worst fit in terms of **decreasing time** and **storage utilization**. Neither first fit nor best fit is clearly better than the other in terms of storage utilization, but first fit is generally faster.

Operating System Concepts – 10th Edition

9.16



Fragmentation 碎片化 (指内存空间的浪费)

- 外部** External Fragmentation – total memory space available to satisfy a request, but it is not contiguous – scattered holes 分散的空洞
 - The storage is fragmented into a large number of small holes.
- 内部** Internal Fragmentation – memory allocated to a process may be larger than requested memory; this size difference is memory internal to a partition, but not being used

Operating System Concepts – 10th Edition

9.17



Fragmentation (Cont.)

- Reduce external fragmentation by a technique called **compaction**
 - Shuffle memory contents to place all free memory together in one large block
- Compaction is possible only if relocation is dynamic and is done at execution time. In another word, if relocation is static and is done at assembly or load time, compaction cannot be done
- The compaction is expensive (time-consuming)
- The backing store (the secondary storage) has similar fragmentation problems, which will be discussed later
- Another more feasible solution is to permit logical address space of the processes to be non-contiguous, by dividing into multiple pieces (variable-sized segments or fixed-size pages), thus allowing a process to be allocated physical memory wherever such a piece of memory is available. These techniques include **segmentation** and **paging**
 - 分段
 - 分页

Operating System Concepts – 10th Edition

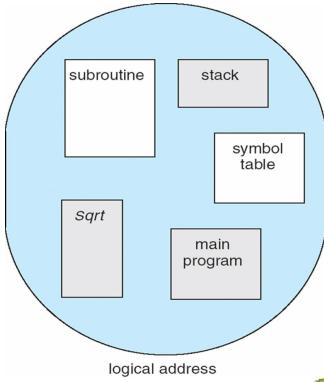
9.18



分段 Segmentation- User View

- Memory-management scheme that supports user view of memory
- A program is a collection of segments. A segment is a logical unit such as:

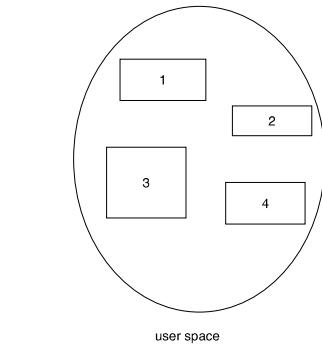
main program
procedure
function
method
object
local variables, global variables
common block
stack
symbol table
arrays



Operating System Concepts – 10th Edition

9.19

Segmentation: Logical vs. Physical



Operating System Concepts – 10th Edition

9.20

Segmentation Architecture

- The Logical address now consists of a **two-tuple**: $\langle \text{segment-number}, \text{offset} \rangle$
- Segment table** – maps two-dimensional programmer-defined addresses (virtual address) into one-dimensional physical addresses; each entry in the segmentation table has:
 - base** – contains the starting physical address where the segments reside in memory
 - limit** – specifies the length of the segment
- Both STBR and STLR stored in the PCB of a process – Recall a process PCB contains all the information about the process
- Segment-table base register (STBR)** points to the segment table's location in memory
- Segment-table length register (STLR)** indicates the number of segments used by a program;

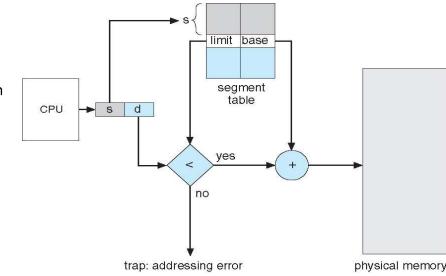
segment number s is legal if $s < \text{STLR}$

Operating System Concepts – 10th Edition

9.21

Segmentation Architecture (Cont.)

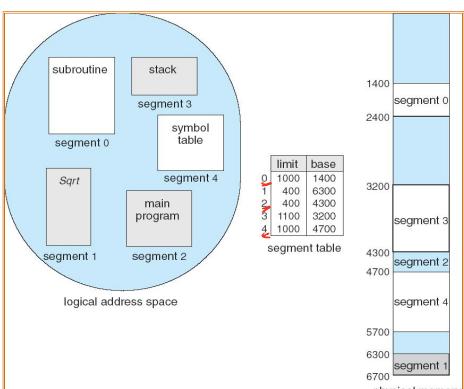
- Protection**, each entry in a segment table associates with:
 - validation bit = 0 \Rightarrow illegal segment
 - read/write/execute privileges
- A protection bit is associated with each segment; **code sharing occurs at segment level**
- Since segments vary in length, memory allocation for a segment is also a **dynamic storage-allocation problem**
- External fragmentation still exists, but is much less severe than that in a contiguous allocation – since each segment is considerably smaller than the total memory space of a process



Operating System Concepts – 10th Edition

9.22

Example of Segmentation



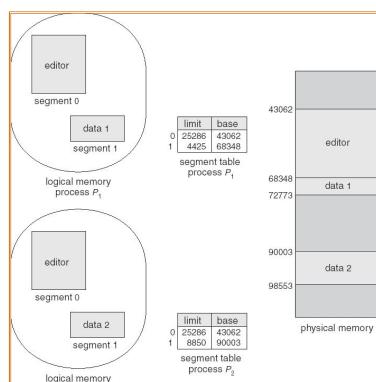
Operating System Concepts – 10th Edition

9.23

Sharing of Segments → 多个进程可以访问同一段代码或数据，通常只发生在只读段中



Sharing of Segments



Operating System Concepts – 10th Edition

9.24

Summary - Segmentation

- Protection is easy in a segmentation scheme – in segment table
 - Code segment would be **read-only**, data and stack would be **read-write** (stores allowed), shared segment could be **read-only** or **read-write**
- Can address be outside the valid range?
 - Yes, this is how stack and heap are allowed to grow. For instance, stack takes fault, system automatically increases the size of stack
- The main problem with segmentation
 - It must fit variable-sized chunks into physical memory – classical **Dynamic Storage Allocation Problem**
 - It may move processes multiple times to fit everything
 - External fragmentation

Operating System Concepts – 10th Edition

9.25



Paging 分页

- Physical address space of a **process** can be **non-contiguous**; process is allocated physical memory whenever the physical memory space is available
- A **frame** is a fixed-size block of memory
- Divide **physical memory** into fixed-sized blocks called **frames**
 - Usually, the size is **power of 2**, between **4KB** (12 bit offset in physical memory address) and **1 GB** (30 bit offset)
- Divide **logical memory** into blocks of **same size** called **pages** – the same size of a frame
- The OS needs to keep track of all free (physical) frames available in main memory
- To run a program of size **N** pages, need to find **N** free frames (non-contiguous) in memory and load the program into that
- A **page table** is used to translate logical to physical addresses – keep track of allocated frames
- This, however, suffers from **internal fragmentation** in the last page/frame

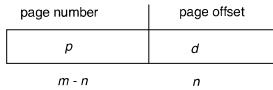
Operating System Concepts – 10th Edition

9.26



Address Translation Scheme

- Address generated by CPU is divided into:
 - **Page number (p)** – used as an index into a **page table** which contains base address of each page in physical memory
 - **Page offset (d)** – combined with base address to define the physical memory address that is sent to the memory unit



- Total m bits logical address, the logical address space 2^m bytes, and the page size or frame size is 2^n , the number of pages this logical address contains is 2^{m-n}

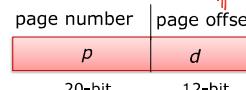
Operating System Concepts – 10th Edition

9.27



Example of Paging Scheme

- Suppose logical address of a system is
 - 1 byte = 8 bits
 - 1 KB = 1024 bytes
 - $2^{12} = 4 \text{ KB}$
- The page size - 12 bits are used to offset into a page, the page size is $2^{12} = 4 \text{ KB}$
- 32 bits virtual address indicates virtual address space is $2^{32} = 4 \text{ GB}$
- 20 bits page number, the number of page is 2^{20}



Operating System Concepts – 10th Edition

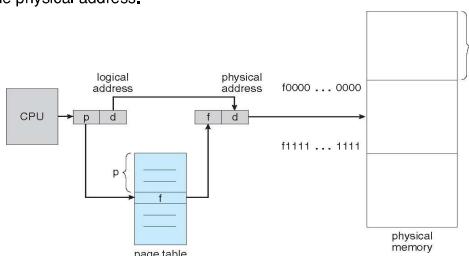
9.28



Paging Hardware

- The following outlines the steps taken by the MMU to translate a logical address generated by the CPU to a physical address:

1. Extract the page number **p** and use it as an index into the page table.
 2. Extract the corresponding frame number **f** from the page table.
 3. Replace the page number **p** in the logical address with the frame number **f**
- As the offset **d** does not change, it is not replaced, and the frame number and offset now comprise the physical address.

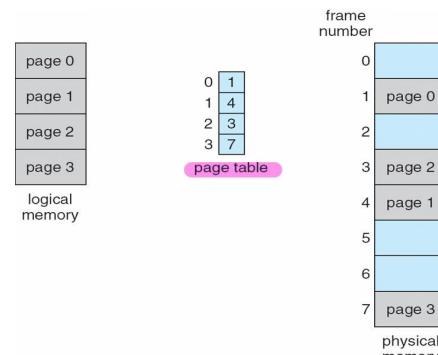


Operating System Concepts – 10th Edition

9.29



Paging Model of Logical and Physical Memory



Operating System Concepts – 10th Edition

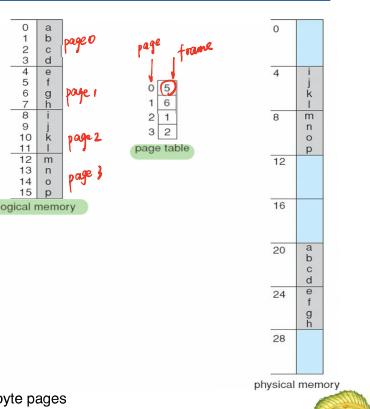
9.30



Paging Example

- Frame number can automatically derive the starting address of each frame $2^4 = 16$ bytes
- $m=4$ represents logical address of 16 bytes
- Logical address 0 is page 0 offset 0. Indexing into the page table, we find that page 0 is in frame 5. Thus, logical address 0 maps to physical address 20 [= $(5 \times 4) + 0$]
- Logical address 3 (page 0, offset 3) maps to physical address 23 [= $(5 \times 4) + 3$].
- Logical address 4 is page 1, offset 0; according to the page table, page 1 is mapped to frame 6. Thus, logical address 4 maps to physical address 24 [= $(6 \times 4) + 0$].
- Logical address 13 maps to physical address 9.

$n=2$ and $m=4$ 32-byte memory and 4-byte pages



Operating System Concepts – 10th Edition

9.31

Paging (Cont.)

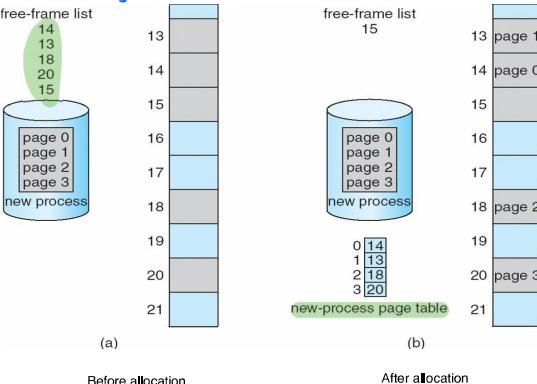
- Calculating internal fragmentation (计算内部碎片)
 - Page size = 2,048 bytes (2KB)
 - Process size = 72,766 bytes
 - 35 pages + 1,086 bytes
 - Internal fragmentation of $2,048 - 1,086 = 962$ bytes
 - Worst case fragmentation = 1 frame – 1 byte
 - On average fragmentation = $1 / 2$ frame size
- So small frame sizes desirable? - not really 框架尺寸不是越小越好！
 - The overhead involved in each page table entry reduces when page size increases.
 - Each page table entry takes memory to track, smaller page size results increased number of page table entries (one per page), thus it leads to larger page table
- Page size has grown over the time
 - SunMicro Solaris OS supports two-page sizes – 8 KB and 4 MB
- Process view and physical memory now very different
 - The programmer views memory as one single space, containing only this one program. But in fact, user program is scattered throughout physical memory, which also holds other programs.

Operating System Concepts – 10th Edition

9.32

Free Frames 空闲帧

在计算机中未被分配给任何进程的内存块



Before allocation

After allocation

Operating System Concepts – 10th Edition

9.33

Implementation of Page Table 页表的实现

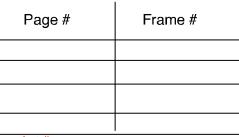
- Page table is kept in main memory – PCB has to keep track of memory allocation for the process
- Page-table base register (PTBR) points to the page table (starting address)
- Page-table length register (PTLR) indicates size of the page table
- In this scheme every data/instruction access requires two memory accesses
 - One for the page table (translation) and one for fetching the data / instruction – same with segmentation scheme 分割方案
- The two memory access problem can be solved by the use of a special fast-lookup hardware cache called associative memory or translation look-aside buffers (TLBs) 联想内存
- TLBs typically small (64 to 1,024 entries). Some CPUs implement separate instruction and data address TLBs. This is shared and used by all processes (kernel and user processes) in a system
- On a TLB miss (if the page number is not in the TLB), value is loaded into the TLB for faster access next time. This can be done by hardware (MMU) or software (running kernel codes)
 - Replacement policies must be considered – LRU, round-robin or even random are possible
 - Some entries can be wired down for permanent fast access, for example TLB entries for key kernel code for fast access

Operating System Concepts – 10th Edition

9.34

Associative Memory 关联内存

- Associative memory (TLB) – parallel search 并行搜索

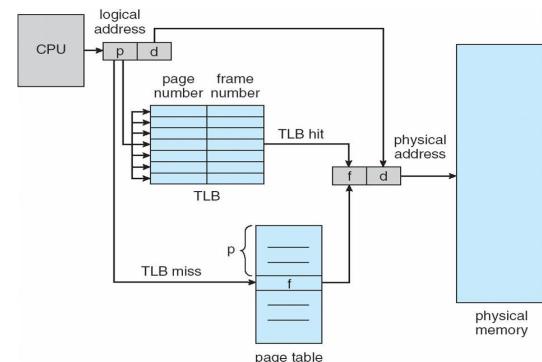


- Address translation (p, d)
 - Check all entries in parallel (hardware) – TLB shared and used by all processes (process ID and page number)
 - If p is in associative register (TLB), get frame # out – TLB hit
 - Otherwise get frame # from page table in memory, and also bring this entry to the TLB
- Locality on TLB
 - Instruction usually stays on the same page (sequential access nature)
 - Stack exhibits locality (pop in and out)
 - Data less locality, still quite a bit

Operating System Concepts – 10th Edition

9.35

Paging Hardware With TLB



Operating System Concepts – 10th Edition

9.36

Effective Access Time

- Associative Lookup = c time unit
 - Usually < 10% of memory access time
- Hit ratio** = α
 - Hit ratio — percentage of times that a page number is found in the **associative registers**; ratio related to number of associative registers
- Effective Access Time (EAT)** (memory access time normalized to 1)

$$EAT = (1 + \epsilon) \alpha + (2 + \epsilon)(1 - \alpha)$$

$$= 2 + \epsilon - \alpha$$
- Consider $\alpha = 80\%$, $\epsilon = 20\text{ns}$ for TLB search, 100ns for memory access
 - EAT = $0.80 \times 120 + 0.20 \times 220 = 140\text{ns}$
- Consider more realistic hit ratio $\rightarrow \alpha = 99\%$, $\epsilon = 20\text{ns}$ for TLB search, 100ns for memory access
 - EAT = $0.99 \times 120 + 0.01 \times 220 = 121\text{ns}$

恐龙

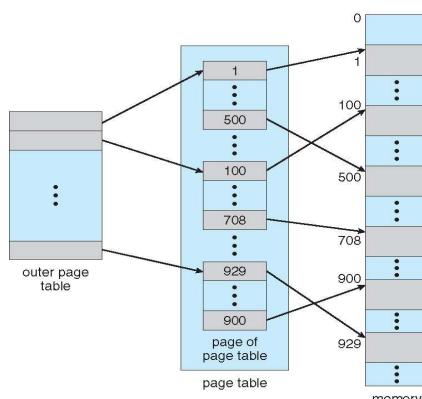


Hierarchical Page Tables 层次页表

- Break up the logical address space into multiple page tables
- A simple technique is a two-level page table or hierarchical page tables
- To page the page table

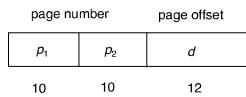


Two-Level Page-Table Scheme



Two-Level Paging Example

- A logical address (on 32-bit machine with 4K page size) is divided into:
 - a page number consisting of 20 bits
 - a page offset consisting of 12 bits
- Since the page table is paged, the page number is further divided into:
 - a 10-bit page number P_2 , **this must be 10 bits if each PTE occupies 4 bytes**, ensuring each inner page table can be stored within one frame (page).
 - a 10-bit page offset P_1 . If PTE also occupies 4 bytes, the number of bits in P_1 **can not be more than 10**, or it needs to be further divided – more than two levels
- Thus, a logical address is as follows:

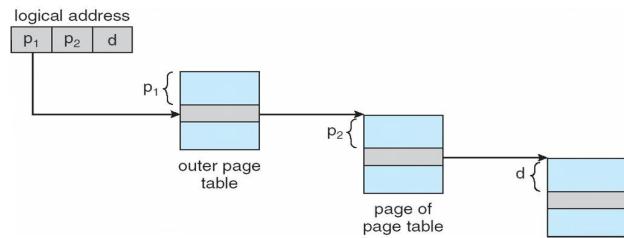


- where p_1 is an index into the **outer page table** (in Intel architecture, this is called "directories", and p_2 is the displacement within the page of the **inner page table**)
- because the address translation works from the outer page table inward, this scheme is also known as a **forward-mapped page table**

正向映射页表

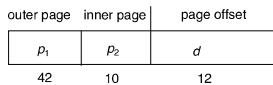


Address-Translation Scheme



64-bit Logical Address Space

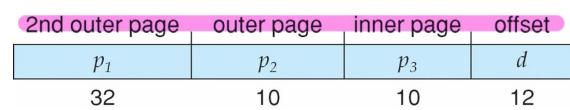
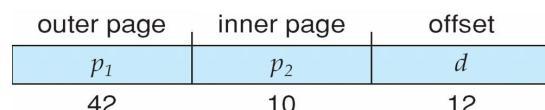
- Even two-level paging scheme not sufficient
- If page size is 4 KB (2^{12})
 - Then page table has 2^{32} entries
 - If two level scheme, inner page tables could be 2^{10} 4-byte entries
 - Address would look like



- Outer page table has 2^{42} entries or 2^{44} bytes
- One solution is to add a 2nd outer page table
- But in the following example the 2nd outer page table is still 2^{34} bytes (16 GB) in size
 - And possibly 4 memory access to get to one physical memory location
 - The 64-bit UltraSPARC would require seven levels of paging – a prohibitive number of memory accesses – to translate each logical address



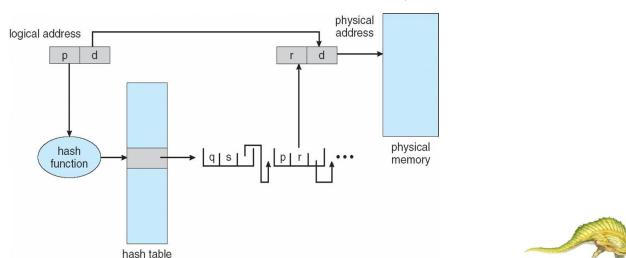
Three-level Paging Scheme





Hashed Page Tables 哈希页表

- Common in address space > 32 bits
- The page number is hashed into a page table
 - This page table contains a chain of elements hashing into the same location
- Each element contains (1) the page number (2) the value of the mapped page frame (3) a pointer to the next element
- Page numbers are compared in this chain searching for a match
 - If a match is found, the corresponding physical frame is extracted

Operating System Concepts – 10th Edition

9.49



Example: The Intel 32 and 64-bit Architectures

- Dominant industry chips 优势行业芯片
 - 16-bit Intel 8086 (late 1970s) and 8088 was used in original IBM PC
- Pentium CPUs are 32-bit and called IA-32 architecture
 - It supports both segmentation and paging - The CPU generates logical addresses, which are given to the segmentation unit. The segmentation unit produces a linear address for each logical address. The linear address is then given to the paging unit, which in turn generates the physical address in main memory.
- Current Intel CPUs are 64-bit and called IA-64 architecture
 - Currently most popular PC operating systems run on Intel chips, including Windows, Mac OS, and Linux (Linux runs on several other architectures as well)
 - Intel's dominance has not spread to mobile systems, where they mainly use ARM architecture
- Many variations in the chips, only the main ideas are covered here

Operating System Concepts – 10th Edition

9.50



Example: The Intel IA-32 Architecture

- Supports both segmentation and segmentation with paging
 - Each segment can be as large as 4GB (32 bits) further divided into pages
 - Each process can have up to 16K segments per process (14 bits), divided into two partitions
 - First partition of up to 8K segments private to process (kept in local descriptor table (LDT))
 - Second partition of up to 8K segments shared among all processes (kept in global descriptor table (GDT))
- CPU generates logical address
 - Selector given to segmentation unit - which produces linear addresses - s or segment number g indicates whether the segment is in LDT or GDT. p deals with protection
- Linear address is then given to the paging unit
 - Which generates physical address in main memory
 - Paging units form equivalent of MMU
 - Pages sizes can be 4 KB or 4 MB

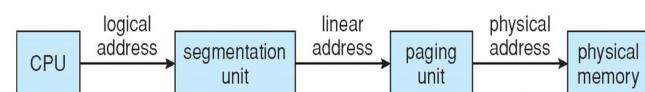
Operating System Concepts – 10th Edition

9.51



Intel IA-32 Segmentation

- Each entry in the LDT and GDT consists of an 8-byte segment descriptor with detailed information about a particular segment, including the base location and limit
- The segment register points to the appropriate entry in the LDT or GDT. The base and limit information about the segment is used to generate a linear address.

Operating System Concepts – 10th Edition

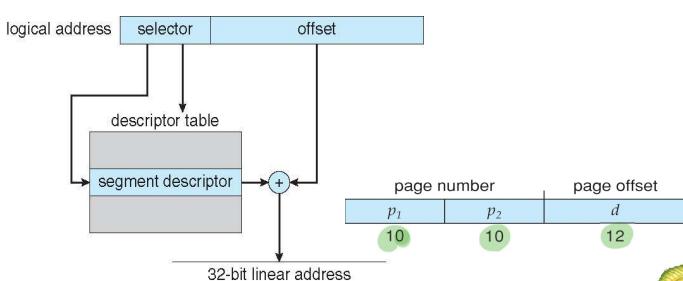
9.52



Intel IA-32 Segmentation

The linear address on the IA-32 is 32 bits long and is formed as follows

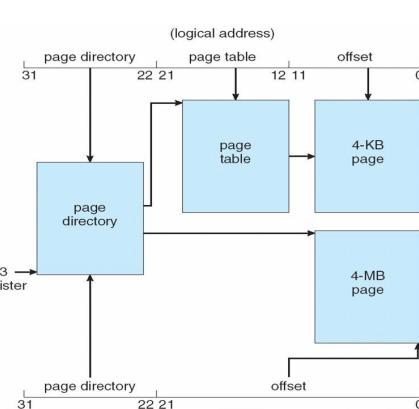
- The limit is used to check for address validity. If the address is not valid, a memory fault is generated, trap to the operating system. If it is valid, the value of the offset is added to the value of the base, resulting in a 32-bit linear address (i.e., each segment can be 4GB)

Operating System Concepts – 10th Edition

9.53



Intel IA-32 Paging Architecture

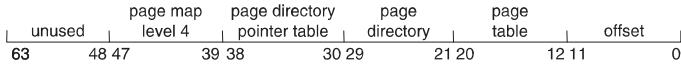
Operating System Concepts – 10th Edition

9.54



Intel x86-64

- Current generation Intel x86-64 architecture
- 64 bits is ginormous (> 16 exabytes)
- In practice only implement 48 bit addressing
 - Page sizes of 4 KB, 2 MB, 1 GB
 - Four levels of paging hierarchy
- Can also use PAE (page address extension) so virtual addresses are 48 bits and physical addresses are 52 bits (4096 terabytes)



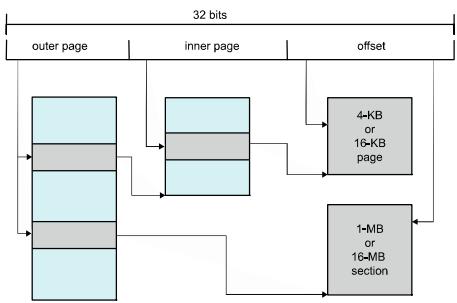
Operating System Concepts – 10th Edition

9.55



Example: ARM Architecture

- Dominant mobile platform chip (Apple iOS and Google Android devices for example)
- Modern, energy efficient, 32-bit CPU
- 4 KB and 16 KB pages
- 1 MB and 16 MB pages (termed sections)
- One-level paging for sections, two-level for smaller pages
- Two levels of TLBs
 - Outer level has two micro TLBs (one data, one instruction)
 - Inner is single main TLB
- First inner is checked, on missouters are checked, and on miss page table walk performed by CPU

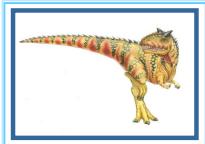


Operating System Concepts – 10th Edition

9.56



End of Chapter 9



Operating System Concepts – 10th Edition