


## TCCS 422: OPERATING SYSTEMS

Memory Virtualization,  
Segmentation,  
Memory Paging

Wes J. Lloyd  
 School of Engineering and Technology,  
 University of Washington - Tacoma



November 26, 2018      TCCS422: Operating Systems [Fall 2018]      Tacoma

School of Engineering and Technology, University of Washington

---

---

---

---

---

---

---

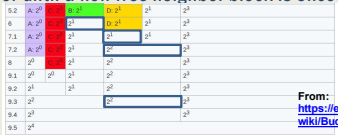
---

---

---

## FEEDBACK FROM 11/20

- Binary buddy allocation:
  - **How does buddy allocation coalesce fragmented memory?**
- **Buddy allocation: Freeing memory blocks:**
  - 1- Free the block of memory
  - 2- Check the neighboring block - is it free too?
  - 3- If free, combine the two, and repeat step 2 until all memory is freed, or until a non-free neighbor block is encountered



From: [https://en.wikipedia.org/wiki/Buddy\\_memory\\_allocation](https://en.wikipedia.org/wiki/Buddy_memory_allocation)

November 26, 2018      TCCS422: Operating Systems [Fall 2018]      L16.2

School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---

---

---

---

---

---

---

## FEEDBACK - 2

- **Which (free space) memory allocation strategy does Ubuntu use?**

Overview from:

- [https://en.wikibooks.org/wiki/The\\_Linux\\_Kernel/Memory](https://en.wikibooks.org/wiki/The_Linux_Kernel/Memory)
- <https://zgallen.github.io/2017/08/03/linux-glibc-mm-overview/>

November 26, 2018      TCCS422: Operating Systems [Fall 2018]      L16.3

School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---

---

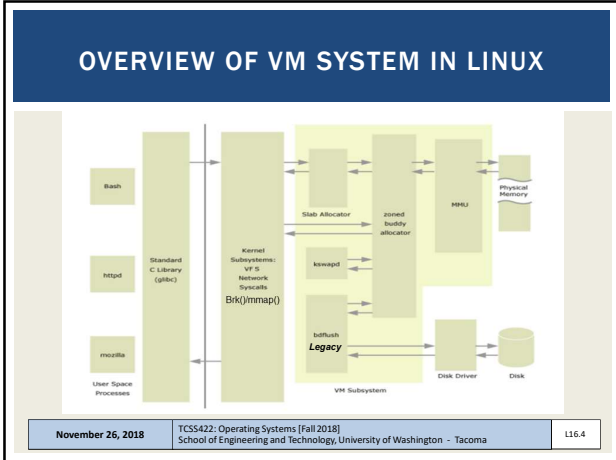
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

### COMPONENTS

- **Memory Management Unit (MMU)** – HW module on CPU, integrates “TLB”, supports virtual memory address translation
- **Buddy Allocator** – Algorithm to allocate/reclaim page frames from physical memory
  - Provides memory pages to consumers such as OS slab allocators (obj caches), kmalloc
  - Page frames managed in a group for buddy allocation in sizes of 2<sup>n</sup> where (size=1,2,4,8,16,32,64,128,256,512,1024...)
  - Memory Zones: DMA/DMA32 (Direct Memory Access) for device I/O, NORMAL, and HIGHMEM (32-bit machines)
  - See /proc/zoneinfo
- **Slab Allocator** – allocates OS object caches – OS structs less than 4kb – provides efficient memory mgmt. for frequently used OS structs

November 26, 2018
L16.5
TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---

---

---

---

---

---

---

### COMPONENTS - 2

- **kswapd** – kernel swap daemon - maintains memory swap space in response to memory demands exceeding physical memory capacity
- Pages can be swapped to disk to reclaim physical memory
- Page frames carry state info to track what to do w/ a page
  - FREE: available
  - ACTIVE: can't swap
  - INACTIVE DIRTY: no longer used, but modified page
  - INACTIVE LAUNDERED: modified page, currently updating to disk
  - INACTIVE CLEAN: no longer being used, can be swapped out
- **Bdflush** – legacy, simple kernel daemon (pdflush thread) to ensure that dirty pages were periodically written to the underlying storage device – now a separate thread is maintained per device

November 26, 2018
L16.6
TCCS422: Operating Systems [Fall 2018]  
 School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---

---

---

---

---

---

---

### PAGING TO DISK ?

- Looking for free space?
- ***What is a likely order of preferred states for selecting a page frame?***
- **Page frame state**
  - FREE: available
  - ACTIVE: can't swap
  - INACTIVE DIRTY: no longer used, but modified page
  - INACTIVE LAUNDERED: modified page, currently updating to disk
  - INACTIVE CLEAN: no longer being used, can be swapped out

November 26, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L16.7
-------------------	---	-------

---

---

---

---

---

---

---

---

### PAGE TRANSLATION EXAMPLE

- ***Can you go over an example of the page table (address) translation?***
- REVIEW Chapter 18...

November 26, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L16.8
-------------------	---	-------

---

---

---

---


---

---

---

---

## REVIEW OF CHAPTER 18: INTRODUCTION TO PAGING



November 26, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L16.9
-------------------	---	-------

---

---

---

---

---

---

---

---

## PAGING

- Split up address space of process into *fixed sized pieces* called **pages**
- Alternative to *variable sized pieces* (Segmentation) which suffers from significant fragmentation
- Physical memory is split up into an array of fixed-size slots called **page frames**.
- Each process has a **page table** which translates virtual addresses to physical addresses

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.10

---

---

---

---

---

---

---

---

---

---

## ADVANTAGES OF PAGING

- Flexibility
  - Abstracts the process address space into pages
  - No need to track direction of HEAP / STACK growth
    - Just add more pages...
  - No need to store unused space
    - As with segments...
- Simplicity
  - Pages and page frames are the same size
  - Easy to allocate and keep a free list of pages

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.11

---

---

---

---

---

---

---

---

---

---

## PAGING: EXAMPLE

- Consider a 128 byte address space with 16-byte pages
- Consider a 64-byte program address space

**Page Table:**  
 VP0 → PF3  
 VP1 → PF7  
 VP2 → PF5  
 VP3 → PF2

A Simple 64-byte Address Space      64-Byte Address Space Placed In Physical Memory

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.12

---

---

---

---

---

---

---

---

---

---

### PAGING: ADDRESS TRANSLATION

- PAGE: Has two address components
  - VPN: Virtual Page Number
  - Offset: Offset within a Page

VPN			offset			
Va5	Va4	Va3	Va2	Va1	Va0	

- Example:  
 Page Size: 16-bytes, Address Space: 64-bytes

VPN			offset			
0	1	0	1	0	1	

Here there are just four pages...

November 26, 2018
L16.13

---

---

---

---

---

---

---

---

---

---

---

---

### EXAMPLE: PAGING ADDRESS TRANSLATION

- Consider a 64-byte program address space (4 pages)
- Stored in 128-byte physical memory (8 frames)
- Offset is preserved
- VPN is looked up

**Page Table:**  
 VP0 → PF3  
 VP1 → PF7  
 VP2 → PF5  
 VP3 → PF2

November 26, 2018
L16.14

---

---

---

---

---

---

---

---

---

---

---

---

### PAGE TRANSLATION EXAMPLE

- Can you go over an example of the page table (address) translation?
- Example:
  - Consider a 64kb computer with 256-byte pages
  - Consider a simple hello world program
    - Program has only 4 memory pages
    - 1 code page, 1 stack page, 1 heap page, 1 data segment page
- (1) How many 256-byte memory pages can the computer hold?
- (VPN) The operating system provides each user program a 64kb virtual address space.
- (2) How many VPN bits are required to index any virtual page?

November 26, 2018
L16.15

---

---

---

---

---

---

---

---

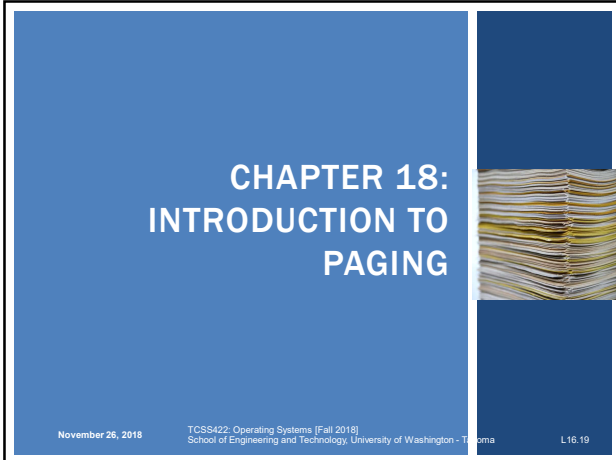
---

---

---

---





CHAPTER 18:  
INTRODUCTION TO  
PAGING

November 26, 2018 TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.19

---

---

---

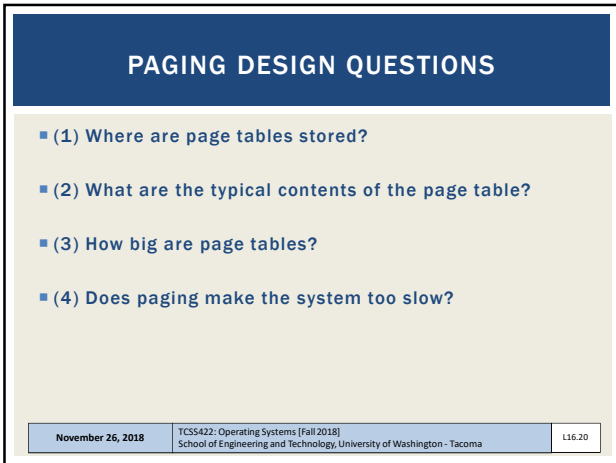
---

---

---

---

---



PAGING DESIGN QUESTIONS

- (1) Where are page tables stored?
- (2) What are the typical contents of the page table?
- (3) How big are page tables?
- (4) Does paging make the system too slow?

November 26, 2018 TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.20

---

---

---

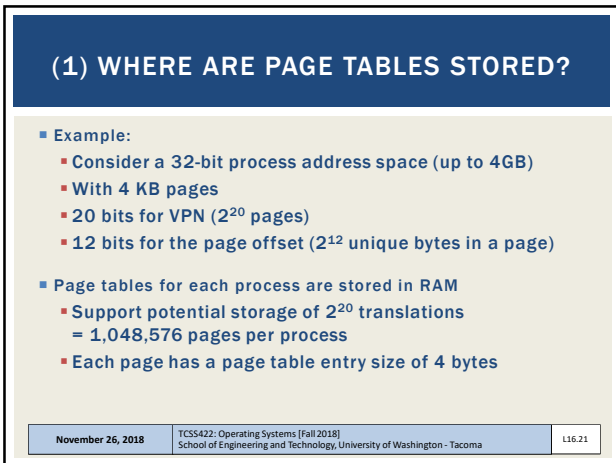
---

---

---

---

---



(1) WHERE ARE PAGE TABLES STORED?

- Example:
  - Consider a 32-bit process address space (up to 4GB)
  - With 4 KB pages
  - 20 bits for VPN ( $2^{20}$  pages)
  - 12 bits for the page offset ( $2^{12}$  unique bytes in a page)
- Page tables for each process are stored in RAM
  - Support potential storage of  $2^{20}$  translations = 1,048,576 pages per process
  - Each page has a page table entry size of 4 bytes

November 26, 2018 TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.21

---

---

---

---

---

---

---

---







### (4) DOES PAGING MAKE THE SYSTEM TOO SLOW?

- Translation
- **Issue #1:** Starting location of the page table is needed
  - HW Support: Page-table base register
    - stores active process
    - Facilitates translation
- **Issue #2:** Each memory address translation for paging requires an extra memory reference
  - HW Support: TLBs (Chapter 19)

**Page Table:**  
 VP0 → PF3  
 VP1 → PF7  
 VP2 → PF5  
 VP3 → PF2

Stored in RAM →

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.28

---

---

---

---

---

---

---

---

---

---

### PAGING MEMORY ACCESS

```

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)
    
```

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.29

---

---

---

---

---

---

---

---

---

---

### COUNTING MEMORY ACCESSES

- Example: Use this Array initialization Code

```
int array[1000];
...
for (i = 0; i < 1000; i++)
    array[i] = 0;
```

- Assembly equivalent:

```
0x1024 movl $0x0, (%edi, %eax, 4)
0x1028 incl %eax
0x102c cmpl $0x03e8, %eax
0x1030 jne 0x1024
```

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.30

---

---

---

---

---

---

---

---

---

---



## OBJECTIVES

- Chapter 19
  - TLB Algorithm
  - TLB Tradeoffs
  - TLB Context Switch

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.34

---

---

---

---

---

---

---

---

## TRANSLATION LOOKASIDE BUFFER

- Legacy name...
- Better name, "Address Translation Cache"
- TLB is an on CPU cache of address translations
  - virtual → physical memory

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.35

---

---

---

---

---

---

---

---

## TRANSLATION LOOKASIDE BUFFER - 2

- Goal:  
Reduce access to the page tables
- Example:  
50 RAM accesses for first 5 for-loop iterations
- Move lookups from RAM to TLB by caching page table entries

The figure consists of three vertically stacked graphs sharing a common x-axis labeled 'Memory Access' ranging from 0 to 50. The top graph shows 'Page Table[39]' (y-axis 1024-1224) and 'Page Table[1]' (y-axis 1024-1074) with scattered access points. The middle graph shows 'Array[VA]' (y-axis 40000-40100) with a dense cluster of accesses between 0-10. The bottom graph shows 'Code[VA]' (y-axis 1024-1124) with a dense cluster of accesses between 0-10.

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.36

---

---

---

---

---

---

---


---







**CHAPTER 20:  
PAGING:  
SMALLER TABLES**



November 26, 2018 TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.46

---

---

---

---

---

---

---

---

**OBJECTIVES**

- Chapter 20
  - Smaller tables
  - Hybrid tables
  - Multi-level page tables

November 26, 2018 TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.47

---

---

---

---

---

---

---

---

**LINEAR PAGE TABLES**

- Consider array-based page tables:
  - Each process has its own page table
  - 32-bit process address space (up to 4GB)
  - With 4 KB pages
  - 20 bits for VPN
  - 12 bits for the page offset

November 26, 2018 TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.48

---

---

---

---

---

---

---

---



### LINEAR PAGE TABLES - 2

- Page tables stored in RAM
- Support potential storage of  $2^{20}$  translations  
= 1,048,576 pages per process @ 4 bytes/page
- Page table size 4MB / process

$$\text{Page table size} = \frac{2^{32}}{2^{12}} * 4\text{Byte} = 4\text{MByte}$$

- Consider 100+ OS processes
  - Requires 400+ MB of RAM to store process information

November 26, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L16.49
-------------------	---	--------

---

---

---

---

---

---

---

---

---

---

### LINEAR PAGE TABLES - 2

- Page tables stored in RAM
- Support potential storage of  $2^{20}$  translations  
= 1,048,576 pages per process @ 4 bytes/page
- Page table size 4MB / process

Page tables are too big and  
consume too much memory.  
Need Solutions ...

- Consider 100+ OS processes
  - Requires 400+ MB of RAM to store process information

November 26, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L16.50
-------------------	---	--------

---

---

---

---

---

---

---

---

---

---

### PAGING: USE LARGER PAGES

- Larger pages** = 16KB =  $2^{14}$
- 32-bit address space:  $2^{32}$
- $2^{18}$  = 262,144 pages

$$\frac{2^{32}}{2^{14}} * 4 = 1\text{MB per page table}$$

- Memory requirement cut to  $\frac{1}{4}$
- However pages are huge
- Internal fragmentation results
- 16KB page(s) allocated for small programs with only a few variables

November 26, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L16.51
-------------------	---	--------

---

---

---

---

---

---

---

---

---

---







### PAGE TABLE INDEX

- 4 bits page directory index (PDI – 1<sup>st</sup> level)
- 4 bits page table index (PTI – 2<sup>nd</sup> level)

14-bits Virtual address

- To dereference one 64-byte memory page,
  - We need one page directory entry (PDE)
  - One page table Index (PTI) – can address 16 pages

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.61

---

---

---

---

---

---

---

---

---

---

---

---

### EXAMPLE - 3

- **For this example, how much space is required to store as a single-level page table with any number of PTEs?**
- 16KB address space, 64 byte pages
- 256 page frames, 4 byte page size
- 1,024 bytes required (*single level*)
- **How much space is required for a two-level page table with only 4 page table entries (PTEs) ?**
- Page directory = 16 entries x 4 bytes (1 x 64 byte page)
- Page table = 4 entries x 4 bytes (1 x 64 byte page)
- 128 bytes required (2 x 64 byte pages)
  - Savings = using just 12.5% the space !!!

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.62

---

---

---

---

---

---

---

---

---

---

---

---

### 32-BIT EXAMPLE

- Consider: 32-bit address space, 4KB pages, 2<sup>20</sup> pages
- Only 4 mapped pages
- **Single level:** 4 MB (we've done this before)
- **Two level:** (old VPN was 20 bits, split in half)
  - Page directory = 2<sup>10</sup> entries x 4 bytes = 1 x 4 KB page
  - Page table = 4 entries x 4 bytes (mapped to 1 4KB page)
  - 8KB (8,192 bytes) required
  - Savings = using just .78 % the space !!!
- 100 sparse processes now require < 1MB for page tables

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.63

---

---

---

---

---

---

---

---

---

---

---

---





## ADDRESS TRANSLATION CODE

```

// 5-level Linux page table address lookup
//
// Inputs:
// mm_struct - process's memory map struct
// vpage - virtual page address

// Define page struct pointers
pgd_t *pgd;
p4d_t *p4d;
pud_t *pud;
pmd_t *pmd;
pte_t *pte;
struct page *page;
    
```

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.70

---

---

---

---

---

---

---

---

---

---

## ADDRESS TRANSLATION - 2

```

pgd = pgd_offset(mm, vpage);
if (pgd_none(*pgd) || pgd_bad(*pgd))
    return 0;
p4d = p4d_offset(pgd, vpage);
if (p4d_none(*p4d) || p4d_bad(*p4d))
    return 0;
pud = pud_offset(p4d, vpage);
if (pud_none(*pud) || pud_bad(*pud))
    return 0;
pmd = pmd_offset(pud, vpage);
if (pmd_none(*pmd) || pmd_bad(*pmd))
    return 0;
if (!(pte = pte_offset_map(pmd, vpage)))
    return 0;
if (!(page = pte_page(*pte)))
    return 0;
physical_page_addr = page_to_phys(page);
pte_unmap(pte);
return physical_page_addr; // param to send back
    
```

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.71

**pgd\_offset():**  
Takes a vpage address and the mm\_struct for the process, returns the PGD entry that covers the requested address...

**p4d/pud/pmd\_offset():**  
Takes a vpage address and the pgd/p4d/pud entry and returns the relevant p4d/pud/pmd.

**pte\_unmap()**  
release temporary kernel mapping for the page table entry

---

---

---

---

---

---


---

---

---

---

## INVERTED PAGE TABLES



- Keep a single page table for each physical page of memory
- Consider 4GB physical memory
- Using 4KB pages, page table requires 4MB to map all of RAM
- Page table stores
  - Which process uses each page
  - Which process virtual page (from process virtual address space) maps to the physical page
- All processes share the same page table for memory mapping, kernel must isolate all use of the shared structure
- Finding process memory pages requires search of 2<sup>20</sup> pages
- Hash table: can index memory and speed lookups

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.72

---

---

---

---

---

---

---

---

---

---



**MULTI-LEVEL PAGE TABLE EXAMPLE**

- Consider a 16 MB computer which indexes memory using 4KB pages
- (#1) For a single level page table, how many pages are required to index memory?
- (#2) How many bits are required for the VPN?
- (#3) Assuming each page table entry (PTE) can index any byte on a 4KB page, how many offset bits are required?
- (#4) Assuming there are 8 status bits, how many bytes are required for each page table entry?

November 26, 2018 TCCS422: Operating Systems [Fall 2018] L16.73  
School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---

---

---

---

---

---

---

**MULTI LEVEL PAGE TABLE EXAMPLE - 2**

- (#5) How many bytes (or KB) are required for a single level page table?
- Let's assume a simple HelloWorld.c program.
- HelloWorld.c requires virtual address translation for 4 pages:
  - 1 - code page                    1 - stack page
  - 1 - heap page                    1 - data segment page
- (#6) Assuming a two-level page table scheme, how many bits are required for the Page Directory Index (PDI)?
- (#7) How many bits are required for the Page Table Index (PTI)?

November 26, 2018 TCCS422: Operating Systems [Fall 2018] L16.74  
School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---

---

---

---

---

---

---

**MULTI LEVEL PAGE TABLE EXAMPLE - 3**

- Assume each page directory entry (PDE) and page table entry (PTE) requires 4 bytes:
  - 6 bits for the Page Directory Index (PDI)
  - 6 bits for the Page Table Index (PTI)
  - 12 offset bits
  - 8 status bits
- (#8) How much **total** memory is required to index the HelloWorld.c program using a two-level page table when we only need to translate 4 total pages?
- **HINT:** we need to allocate one Page Directory and one Page Table...
- **HINT:** how many entries are in the PD and PT

November 26, 2018 TCCS422: Operating Systems [Fall 2018] L16.75  
School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---

---

---

---

---

---

---

### MULTI LEVEL PAGE TABLE EXAMPLE - 4

- (#9) Using a single page directory entry (PDE) pointing to a single page table (PT), if all of the slots of the page table (PT) are in use, what is the total amount of memory a two-level page table scheme can address?
- (#10) And finally, for this example, as a percentage (%), how much memory does the 2-level page table scheme consume compared to the 1-level scheme?
- HINT: two-level memory use / one-level memory use

November 26, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.76

---

---

---

---

---

---

---

---

### ANSWERS

- #1 - 4096 pages
- #2 - 12 bits
- #3 - 12 bits
- #4 - 4 bytes
- #5 -  $4096 \times 4 = 16,384$  bytes (16KB)
- #6 - 6 bits
- #7 - 6 bits
- #8 - 256 bytes for Page Directory (PD) (64 entries x 4 bytes)  
256 bytes for Page Table (PT) **TOTAL = 512 bytes**
- #9 - 64 entries, where each entry maps a 4,096 byte page  
With 12 offset bits, can address 262,144 bytes (256 KB)
- #10-  $512/16384 = .03125 \rightarrow 3.125\%$

November 26, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.77

---

---

---

---


---

---

---

---

### CHAPTER 21/22: BEYOND PHYSICAL MEMORY



November 26, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.78

---

---

---

---

---

---

---

---

## MEMORY HIERARCHY

- Disks (HDD, SSD) provide another level of storage in the memory hierarchy

Memory Hierarchy in modern system

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.79

---

---

---

---

---

---

---

---

---

---

## MOTIVATION FOR EXPANDING THE ADDRESS SPACE

- Can provide illusion of an address space larger than physical RAM
- For a single process
  - Convenience
  - Ease of use
- For multiple processes
  - Large virtual memory space for many concurrent processes

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.80

---

---

---

---

---

---

---

---

---

---

## LATENCY TIMES

- Design considerations
  - SSDs 4x the time of DRAM
  - HDDs 80x the time of DRAM

Action	Latency (ns)	(µs)	
L1 cache reference	0.5ns		
L2 cache reference	7 ns		14x L1 cache
Mutex lock/unlock	25 ns		
Main memory reference	100 ns		20x L2 cache, 200x L1
Read 4K randomly from SSD*	150,000 ns	150 µs	~1GB/sec SSD
Read 1 MB sequentially from memory	250,000 ns	250 µs	
Read 1 MB sequentially from SSD*	1,000,000 ns	1,000 µs	1 ms ~1GB/sec SSD, 4X memory
Read 1 MB sequentially from disk	20,000,000 ns	20,000 µs	20 ms 80x memory, 20X SSD

- Latency numbers every programmer should know
- From: <https://gist.github.com/jboner/2841832#file-latency-txt>

November 26, 2018
TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.81

---

---

---

---

---

---

---

---

---

---

## SWAP SPACE

- Disk space for storing memory pages
- “Swap” them in and out of memory to disk as needed

	PFN 0	PFN 1	PFN 2	PFN 3
<b>Physical Memory</b>	Proc 0 [VPN 0]	Proc 1 [VPN 2]	Proc 1 [VPN 3]	Proc 2 [VPN 0]
<b>Swap Space</b>	Block 0 Proc 0 [VPN 1]	Block 1 Proc 0 [VPN 2]	Block 2 [Free]	Block 3 Proc 1 [VPN 0]
	Block 4 Proc 1 [VPN 1]	Block 5 Proc 3 [VPN 0]	Block 6 Proc 2 [VPN 1]	Block 7 Proc 3 [VPN 1]

Physical Memory and Swap Space

November 26, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.82

---

---

---

---

---

---

---

---

---

---

---

---

## PAGE LOCATION

- Page table pages are:
  - Stored in memory
  - Swapped to disk
- Present bit
  - In the page table entry (PTE) indicates if page is present
- Page fault
  - Memory page is accessed, but has been swapped to disk

November 26, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.83

---

---

---

---

---

---

---

---

---

---

---

---

## PAGE FAULT

- OS steps in to handle the page fault
- Loading page from disk requires a free memory page
- Page-Fault Algorithm

```

1: PFN = FindFreePhysicalPage ()
2: if (PFN == -1) // no free page found
3:     PFN = EvictPage () // run replacement algorithm
4:     DiskRead (PTE.DiskAddr, pfn) // sleep (waiting for I/O)
5:     PTE.present = True // set PTE bit to present
6:     PTE.PFN = PFN // reference new loaded page
7:     RetryInstruction () // retry instruction
            
```

November 26, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.84

---

---

---

---

---

---

---

---

---

---

---

---

## PAGE REPLACEMENTS

- Page daemon
  - Background threads which monitors swapped pages
  
- Low watermark (LW)
  - Threshold for when to swap pages to disk
  - Daemon checks: free pages < LW
  - Begin swapping to disk until reaching the highwater mark
  
- High watermark (HW)
  - Target threshold of free memory pages
  - Daemon free until: free pages >= HW

November 26, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.85

---

---

---

---

---

---


---

---

---

---

## REPLACEMENT POLICIES



November 26, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.86

---

---

---

---

---

---

---

---

---

---

## CACHE MANAGEMENT

- Replacement policies apply to “any” cache
- Goal is to minimize the number of misses
- Average memory access time can be estimated:

$$AMAT = (P_{Hit} * T_M) + (P_{Miss} * T_D)$$

Argument	Meaning
$T_M$	The cost of accessing memory (time)
$T_D$	The cost of accessing disk (time)
$P_{Hit}$	The probability of finding the data item in the cache(a hit)
$P_{Miss}$	The probability of not finding the data in the cache(a miss)

- Consider  $T_M = 100 \text{ ns}$ ,  $T_D = 10 \text{ ms}$
- Consider  $P_{hit} = .9$  (90%),  $P_{miss} = .1$
- Consider  $P_{hit} = .999$  (99.9%),  $P_{miss} = .001$

November 26, 2018
TCSS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma
L16.87

---

---

---

---

---

---

---

---

---

---



### HISTORY-BASED POLICIES

- LRU: Least recently used
  - Always replace page with oldest access time (front)
  - Always move end of cache when element is read again
  - Considers temporal locality (*when pg was last accessed*)

0 1 2 0 1 3 0 3 1 2 1

**What is the hit/miss ratio?**  
**6 hits**

- LFU: Least frequently used
  - Always replace page with fewest accesses (front)
  - Consider frequency of page accesses

0 1 2 0 1 3 0 3 1 2 1

**Hit/miss ratio is =**  
**6 hits**

November 26, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.91

---

---

---

---

---

---

---

---

---

---

---

---

### WORKLOAD EXAMPLES: NO-LOCALITY

- No-Locality (Random Access) Workload
  - Perform 10,000 random page accesses
  - Across set of 100 memory pages

When the cache is large enough to fit the entire workload, it doesn't matter which policy you use.

November 26, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.92

---

---

---

---

---

---

---

---

---

---

---

---

### WORKLOAD EXAMPLES: 80/20

- 80/20 Workload
  - Perform 10,000 page accesses, against set of 100 pages
  - 80% of accesses are to 20% of pages (hot pages)
  - 20% of accesses are to 80% of pages (cold pages)

LRU is more likely to hold onto hot pages  
(recalls history)

November 26, 2018 TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.93

---

---

---

---

---

---

---

---

---

---

---

---

### WORKLOAD EXAMPLES: SEQUENTIAL

- Looping sequential workload
  - Refer to 50 pages in sequence: 0, 1, ..., 49
  - Repeat loop

Random performs better than FIFO and LRU for cache sizes < 50

Algorithms should provide "scan resistance"

November 26, 2018 TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.94

---

---

---

---

---

---

---

---

---

---

### IMPLEMENTING LRU

- Implementing last recently used (LRU) requires tracking access time for all system memory pages
- Times can be tracked with a list
- For cache eviction, we must scan an entire list
- Consider: 4GB memory system ( $2^{32}$ ), with 4KB pages ( $2^{12}$ )
- This requires  $2^{20}$  comparisons !!!
- Simplification is needed
  - Consider how to approximate the oldest page access

November 26, 2018 TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.95

---

---

---

---

---

---

---

---

---

---

### IMPLEMENTING LRU - 2

- Harness the Page Table Entry (PTE) Use Bit
- HW sets to 1 when page is used
- OS sets to 0
- Clock algorithm (approximate LRU)
  - Refer to pages in a circular list
  - Clock hand points to current page
  - Loops around
    - IF USE\_BIT=1 set to USE\_BIT = 0
    - IF USE\_BIT=0 replace page

November 26, 2018 TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma L16.96

---

---

---

---

---

---

---

---

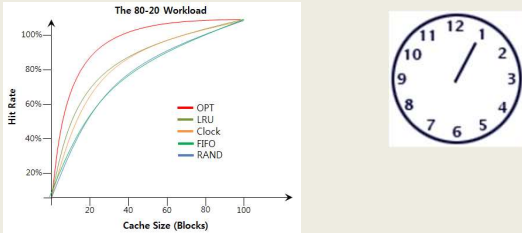
---

---



### CLOCK ALGORITHM

- Not as efficient as LRU, but better than other replacement algorithms that do not consider history



November 26, 2018 TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.97

---

---

---

---

---

---

---

---

---

---

### CLOCK ALGORITHM - 2

- Consider dirty pages in cache
  - If DIRTY (modified) bit is FALSE
    - No cost to evict page from cache
  - If DIRTY (modified) bit is TRUE
    - Cache eviction requires updating memory
    - Contents have changed
- Clock algorithm should favor no cost eviction

November 26, 2018 TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.98

---

---

---

---

---

---

---

---

---

---

### WHEN TO LOAD PAGES

- On demand → demand paging
- Prefetching
  - Preload pages based on anticipated demand
  - Prediction based on locality
  - Access page P, suggest page P+1 may be used
- What other techniques might help anticipate required memory pages?
  - Prediction models, historical analysis
  - In general: accuracy vs. effort tradeoff
  - High analysis techniques struggle to respond in real time

November 26, 2018 TCCS422: Operating Systems [Fall 2018]  
School of Engineering and Technology, University of Washington - Tacoma L16.99

---

---

---

---

---

---

---

---

---

---

**OTHER SWAPPING POLICIES**

- Page swaps / writes
  - Group/cluster pages together
  - Collect pending writes, perform as batch
  - Grouping disk writes helps amortize latency costs
- Thrashing
  - Occurs when system runs many memory intensive processes and is low in memory
  - Everything is constantly swapped to-and-from disk

November 26, 2018 TCCS422: Operating Systems [Fall 2018] L16.100  
School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---

---

---

---

---

**OTHER SWAPPING POLICIES - 2**

- Working sets
  - Groups of related processes
  - When thrashing: prevent one or more working set(s) from running
  - Temporarily reduces memory burden
  - Allows some processes to run, reduces thrashing

November 26, 2018 TCCS422: Operating Systems [Fall 2018] L16.101  
School of Engineering and Technology, University of Washington - Tacoma

---

---

---

---


---

---

---

---

**QUESTIONS**



---

---

---

---

---

---

---

---