# TCSS 422: OPERATING SYSTEMS

**Memory Virtualization,
Segmentation,
Memory Paging**

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

**November 26, 2018**

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington    Tacoma

---

# 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**

**November 26, 2018**

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington  - Tacoma

L16.2

# FEEDBACK - 2

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

- Overview from:
- https://en.wikibooks.org/wiki/The_Linux_Kernel/Memory
- https://zgqallen.github.io/2017/08/03/linux-glic-mm-overview/

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.3 |
|---|---|---|

# OVERVIEW OF VM SYSTEM IN LINUX



| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.4 |
|---|---|---|

## 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 $2^n$ 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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.5 |
|---|---|---|

## 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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.6 |
|---|---|---|

# 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]<br>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…**

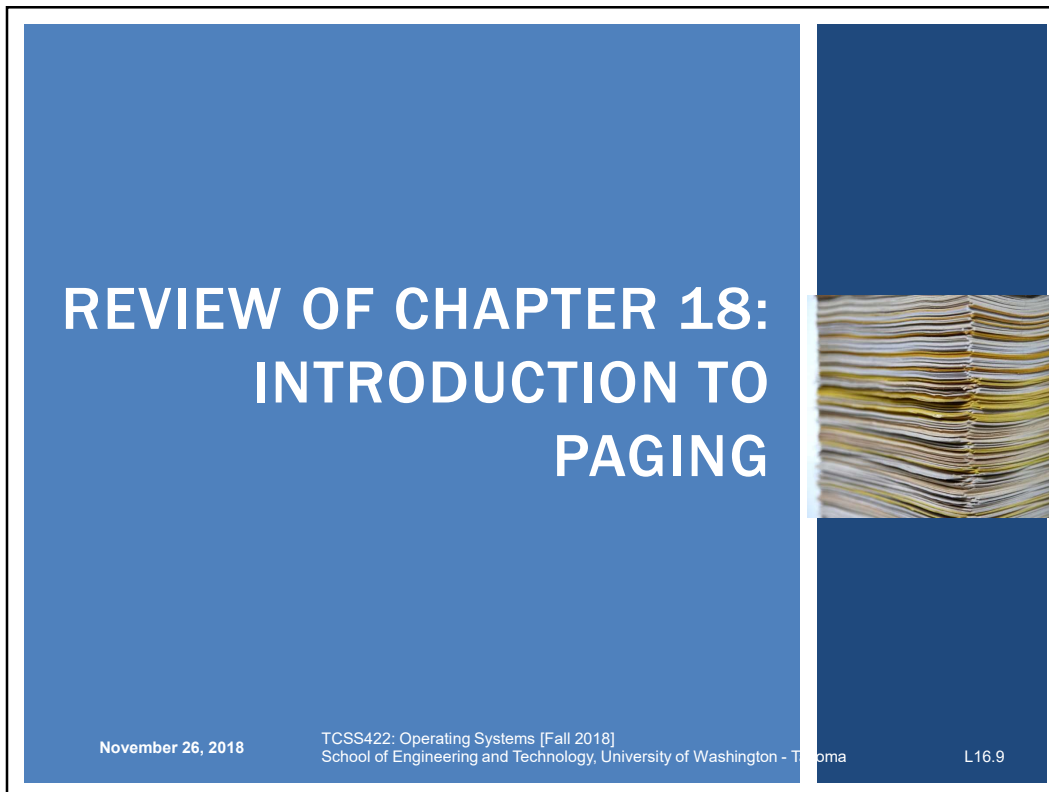| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.8 |
|---|---|---|

# REVIEW OF CHAPTER 18: INTRODUCTION TO PAGING

# 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

## 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
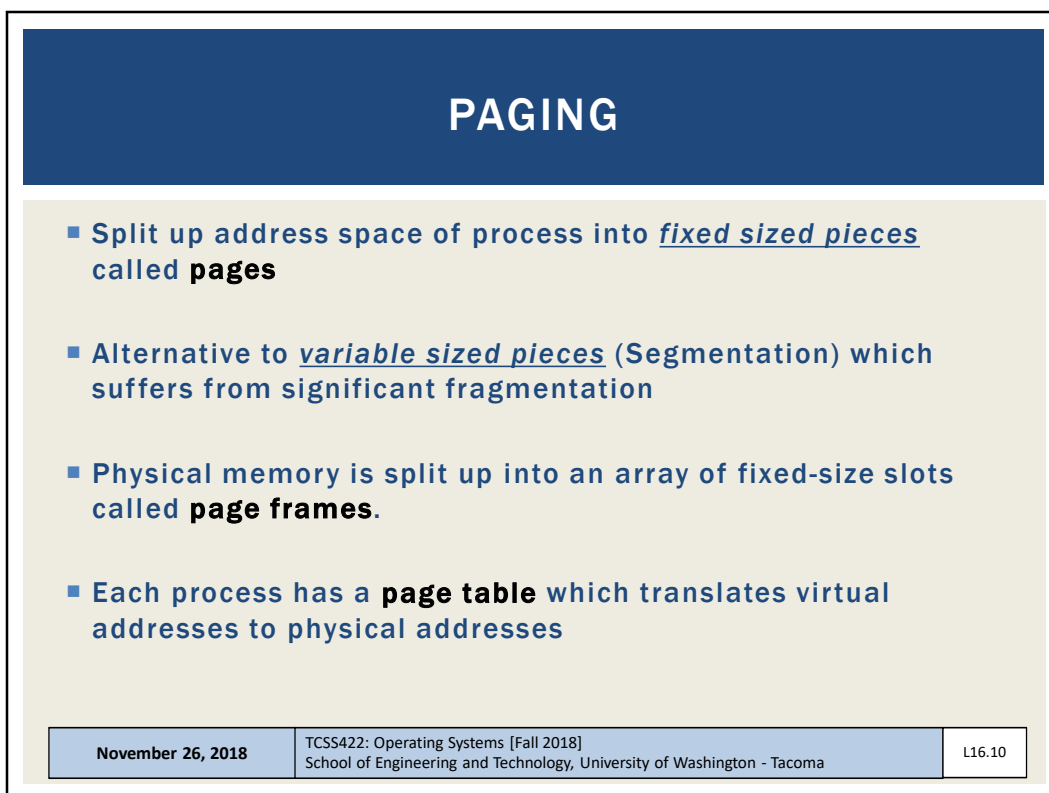  - Easy to allocate and keep a free list of pages

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.11 |

## PAGING: EXAMPLE

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

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

- Consider a 64-byte program address space



A Simple 64-byte Address Space

64-Byte Address Space Placed In Physical Memory

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.13 |
|---|---|---|

## EXAMPLE:<br>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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | 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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.15 |
| --- | --- | --- |

## EXAMPLE - 2

- (3) To reference any individual byte on a 256-byte page, how many bits are required (OFFSET bits)?

- A single-level page table provides a one-dimensional array to look up the physical frame number of each virtual memory page
- Each page table entry (PTE) is like a record. It contains the Physical Frame Number (PFN) and status bits for the page
- PTE example with 20-bit PTE:

| 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|
| PFN | G | PAT | D | A | PCD | PWT | U/S | R/W | P |

An x86 Page Table Entry(PTE)

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.16 |
| --- | --- | --- |

## EXAMPLE - 3

- Now consider our Page Table Entry (PTE) for our 64kb computer

- (4) How bits are required for the PFN?

- (5) Assuming there are 8 status bits, what is the PTE size in bits? Bytes?

- (6) What is the storage requirement for a 1-level page table?

- (7) Using 1-level page tables to index memory, how many process would fill main memory with page tables!!??

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.17 |

## OBJECTIVES

- Quiz 4
- Quiz 5
- Program 3

- **Paging**
- Chapter 18 – Introduction to Paging (*finish…*)
- Chapter 19 – Translation Lookaside Buffer
- Chapter 20 – Paging Smaller Tables
- Chapter 21/22 – Beyond Physical Memory

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.18 |

# CHAPTER 18: INTRODUCTION TO PAGING

# 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?

## (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

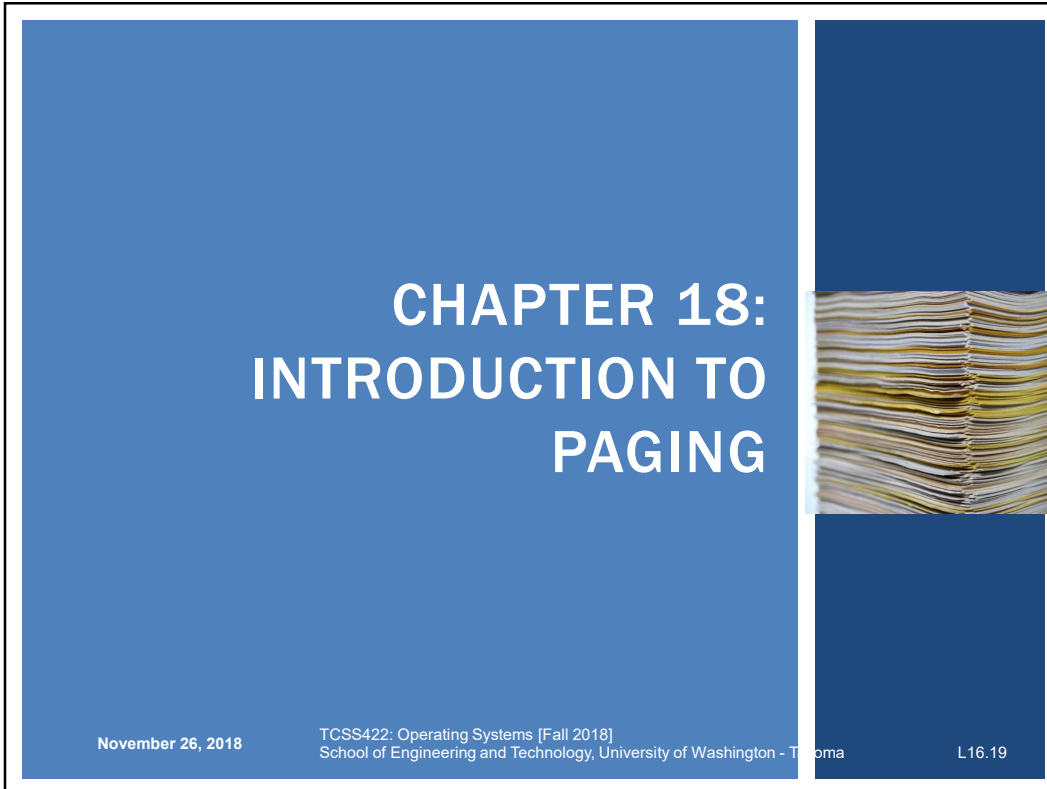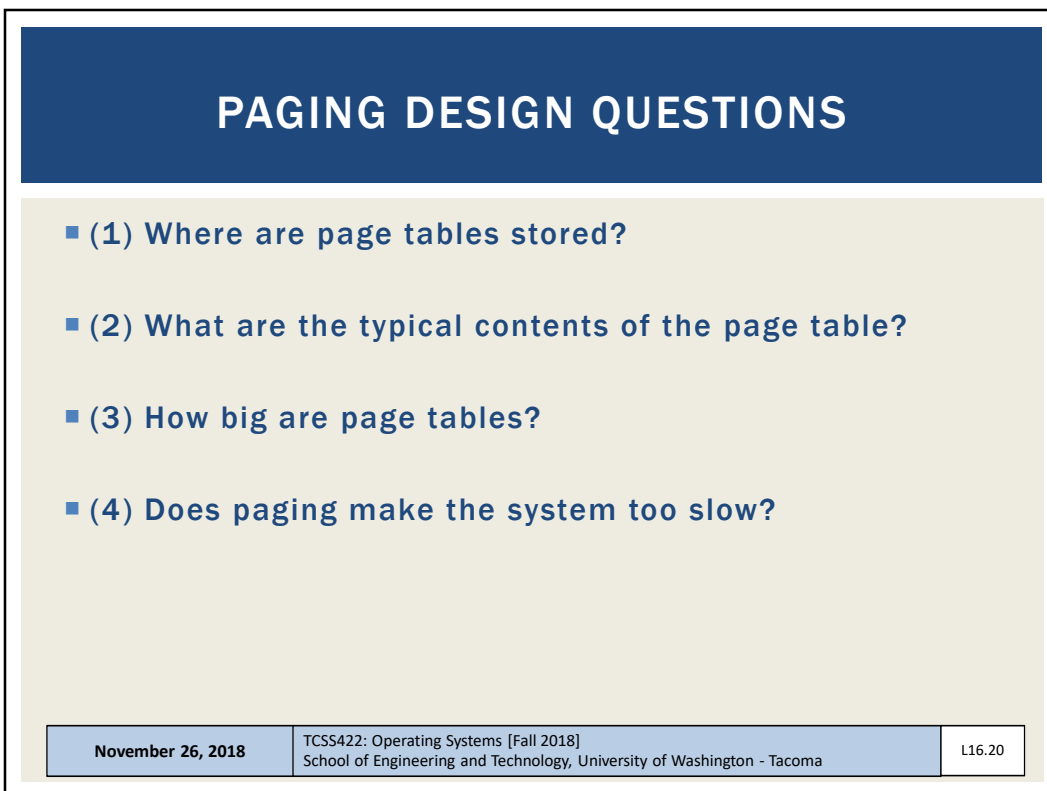| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.21 |
|---|---|---|

## PAGE TABLE EXAMPLE

- With $2^{20}$ slots in our page table for a single process

- Each slot dereferences a VPN

- Provides physical frame number

- Each slot requires 4 bytes (32 bits)
  - 20 for the PFN on a 4GB system with 4KB pages
  - 12 for the offset which is preserved
  - (note we have no status bits, so this is unrealistically small)

| VPN$_0$ |
|---|
| VPN$_1$ |
| VPN$_2$ |
| ... |
| ... |
| VPN$_{1048576}$ |

- How much memory to store page table for 1 process?
  - 4,194,304 bytes (or 4MB) to index one process

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.22 |
|---|---|---|

## NOW FOR AN ENTIRE OS

- If 4 MB is required to store one process

- Consider how much memory is required for an entire OS?
  - With for example 100 processes…

- Page table memory requirement is now 4MB x 100 = 400MB

- If computer has 4GB memory (maximum for 32-bits), the page table consumes 10% of memory

  400 MB / 4000 GB

- Is this efficient?

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.23 |

## (2) WHAT'S ACTUALLY IN THE PAGE TABLE

- Page table is data structure used to map virtual page numbers (VPN) to the physical address (Physical Frame Number PFN)
  - Linear page table → simple array

- Page-table entry
  - 32 bits for capturing state

| 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 | 11 10 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| PFN | | G | PAT | D | A | PCD | PWT | U/S | R/W | P |

An x86 Page Table Entry(PTE)

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.24 |

# PAGE TABLE ENTRY

- P: present
- R/W: read/write bit
- U/S: supervisor
- A: accessed bit
- D: dirty bit
- PFN: the page frame number

| 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 | 11 10 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| PFN | | G | PAT | D | A | PCD | PWT | U/S | R/W | P |

An x86 Page Table Entry(PTE)

# PAGE TABLE ENTRY - 2

- Common flags:

- **Valid Bit:** Indicating whether the particular translation is valid.

- **Protection Bit:** Indicating whether the page could be read from, written to, or executed from

- **Present Bit:** Indicating whether this page is in physical memory or on disk(swapped out)

- **Dirty Bit:** Indicating whether the page has been modified since it was brought into memory

- **Reference Bit(Accessed Bit):** Indicating that a page has been accessed

# (3) HOW BIG ARE PAGE TABLES?

- Page tables are too big to store on the CPU

- Page tables are stored using physical memory

- Paging supports efficiently storing a sparsely populated address space

  - Reduced memory requirement
    Compared to base and bounds, and segments

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.27 |

# (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

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

- **Issue #2:** Each memory address translation for paging requires an extra memory reference
  - HW Support: TLBs (Chapter 19)

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: 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 | TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma | L16.30 |
|---|---|---|

# VISUALIZING MEMORY ACCESSES:
## FOR THE FIRST 5 LOOP ITERATIONS

- Locations:
  - Page table
  - Array
  - Code

- 50 accesses for 5 loop iterations

# PAGING SYSTEM EXAMPLE

- Consider a 4GB Computer:
- With a 4096-byte page size (4KB)
- How many pages would fit in physical memory?

- Now consider a page table:
- For the page table entry, how many bits are required for the VPN?
- If we assume the use of 4-byte (32 bit) page table entries, how many bits are available for status bits?
- How much space does this page table require?
  Page Table Entries x Number of pages
- How many page tables (for user processes) would fill the entire 4GB of memory?

# CHAPTER 19: TRANSLATION LOOKASIDE BUFFER (TLB)

November 26, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L16.33

# OBJECTIVES

- Chapter 19

  - TLB Algorithm

  - TLB Tradeoffs

  - TLB Context Switch

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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

# 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

# TRANSLATION LOOKASIDE BUFFER (TLB)

- Part of the CPU's Memory Management Unit (MMU)

- Address translation cache



**Address Translation with MMU**

# TRANSLATION LOOKASIDE BUFFER (TLB)

- Part of the CPU's Memory Management Unit (MMU)

- Address translation cache



The TLB is an address translation cache
Different than L1, L2, L3 CPU memory caches

**Address Translation with MMU**

# TLB BASIC ALGORITHM

- **For: array based page table**
- **Hardware managed TLB**

```
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:        AccessMemory( PhysAddr )
8:    }else RaiseException(PROTECTION_ERROR)
```

**Generate the physical address to access memory**

# TLB BASIC ALGORITHM - 2

```
11:    else{ //TLB Miss
12:        PTEAddr = PTBR + (VPN * sizeof(PTE))
13:        PTE = AccessMemory(PTEAddr)
14:        (…)  // Check for, and raise exceptions…
15:
16:        TLB_Insert( VPN , PTE.PFN , PTE.ProtectBits)
17:        RetryInstruction()
18:    }
19:}
```

**Retry the instruction... (requery the TLB)**

# TLB – ADDRESS TRANSLATION CACHE

- Key detail:

- For a TLB miss, we first access the page table in RAM to populate the TLB... **we then requery the TLB**
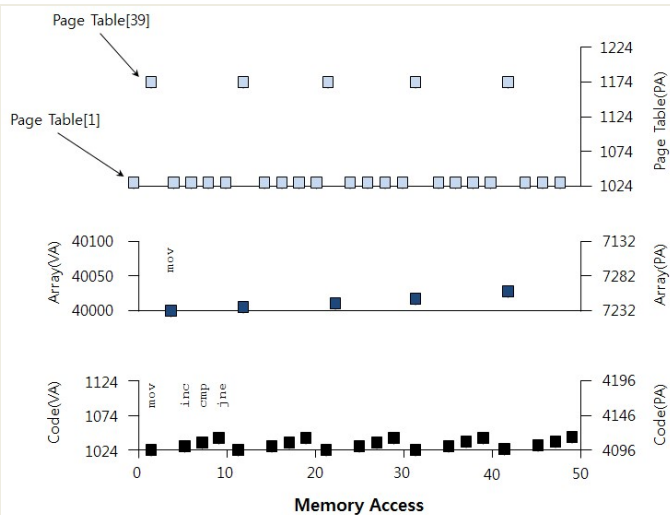
- **All address translations go through the TLB**

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.41 |
|---|---|---|

# TLB EXAMPLE

```
0:          int sum = 0 ;
1:          for( i=0; i<10; i++){
2:                  sum+=a[i];
3:          }
```

- Example:

- Program address space: 256-byte
  - Addressable using 8 total bits $(2^8)$
  - 4 bits for the VPN (16 total pages)

- Page size: 16 bytes
  - Offset is addressable using 4-bits

- Store an array: of (10) 4-byte integers

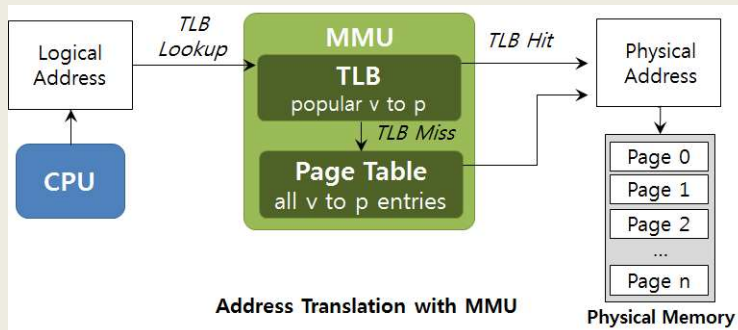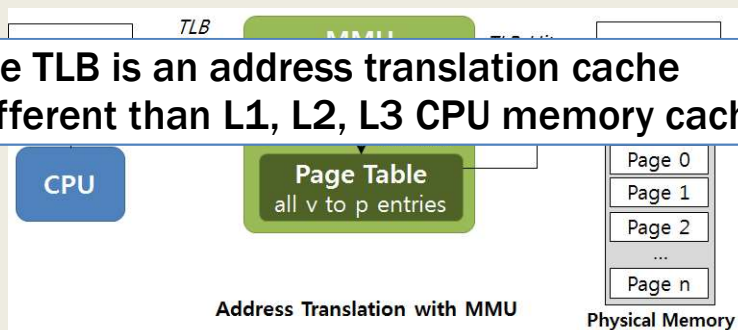| | OFFSET | | | | |
|---|---|---|---|---|---|
| | 00 | 04 | 08 | 12 | 16 |
| VPN = 00 | | | | | |
| VPN = 01 | | | | | |
| VPN = 03 | | | | | |
| VPN = 04 | | | | | |
| VPN = 05 | | | | | |
| VPN = 06 | | a[0] | a[1] | a[2] | |
| VPN = 07 | a[3] | a[4] | a[5] | a[6] | |
| VPN = 08 | a[7] | a[8] | a[9] | | |
| VPN = 09 | | | | | |
| VPN = 10 | | | | | |
| VPN = 11 | | | | | |
| VPN = 12 | | | | | |
| VPN = 13 | | | | | |
| VPN = 14 | | | | | |
| VPN = 15 | | | | | |

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.42 |
|---|---|---|

# TLB EXAMPLE - 2

```
0:      int sum = 0 ;
1:      for( i=0; i<10; i++){
2:          sum+=a[i];
3:      }
```

OFFSET
00   04   08   12   16

VPN = 00
VPN = 01
VPN = 03
VPN = 04
VPN = 05
VPN = 06   | a[0] | a[1] | a[2]
VPN = 07   a[3] | a[4] | a[5] | a[6]
VPN = 08   a[7] | a[8] | a[9]
VPN = 09
VPN = 10
VPN = 11
VPN = 12
VPN = 13
VPN = 14
VPN = 15

- Consider the code above:

- Initially the TLB does not know where a[] is
- Consider the accesses:
- a[0], a[1], a[2], a[3], a[4], a[5], a[6], a[7], a[8], a[9]
- How many pages are accessed?
- What happens when accessing a page not in the TLB?

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.43 |

# TLB EXAMPLE - 3

```
0:      int sum = 0 ;
1:      for( i=0; i<10; i++){
2:          sum+=a[i];
3:      }
```

OFFSET
00   04   08   12   16

VPN = 00
VPN = 01
VPN = 03
VPN = 04
VPN = 05
VPN = 06   | a[0] | a[1] | a[2]
VPN = 07   a[3] | a[4] | a[5] | a[6]
VPN = 08   a[7] | a[8] | a[9]
VPN = 09
VPN = 10
VPN = 11
VPN = 12
VPN = 13
VPN = 14
VPN = 15

- For the accesses: a[0], a[1], a[2], a[3], a[4],
- a[5], a[6], a[7], a[8], a[9]

- How many are hits?
- How many are misses?
- What is the hit rate? (%)
  - 70% (3 misses one for each VP, 7 hits)

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.44 |

## TLB EXAMPLE - 4

```
0:      int sum = 0 ;
1:      for( i=0; i<10; i++){
2:              sum+=a[i];
3:      }
```

- **What factors affect the hit/miss rate?**
  - **Page size**
  - **Data locality**
  - **Temporal locality**

| | OFFSET | | | | |
|---|---|---|---|---|---|
| | 00 | 04 | 08 | 12 | 16 |
| VPN = 00 | | | | | |
| VPN = 01 | | | | | |
| VPN = 03 | | | | | |
| VPN = 04 | | | | | |
| VPN = 05 | | | | | |
| VPN = 06 | | a[0] | a[1] | a[2] | |
| VPN = 07 | a[3] | a[4] | a[5] | a[6] | |
| VPN = 08 | a[7] | a[8] | a[9] | | |
| VPN = 09 | | | | | |
| VPN = 10 | | | | | |
| VPN = 11 | | | | | |
| VPN = 12 | | | | | |
| VPN = 13 | | | | | |
| VPN = 14 | | | | | |
| VPN = 15 | | | | | |

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.45 |
|---|---|---|

# CHAPTER 20:
# PAGING:
# SMALLER TABLES

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

# OBJECTIVES

- Chapter 20

  - Smaller tables

  - Hybrid tables

  - Multi-level page tables

# 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

## 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

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

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

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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 <u>too big</u> and consume too much memory.
>
> Need Solutions ...

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

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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

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

- Memory requirement cut to ¼
- However pages are huge
- Internal fragmentation results
- 16KB page(s) allocated for small programs with only a few variables

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.51 |
|---|---|---|

# PAGE TABLES: WASTED SPACE

- **Process: 16KB Address Space w/ 1KB pages**



A 16KB Address Space with 1KB Pages

| PFN | valid | prot | present | dirty |
|---|---|---|---|---|
| 10 | 1 | r-x | 1 | 0 |
| - | 0 | - | - | - |
| - | 0 | - | - | - |
| - | 0 | - | - | - |
| 15 | 1 | rw- | 1 | 1 |
| ... | ... | ... | ... | ... |
| - | 0 | - | - | - |
| 3 | 1 | rw- | 1 | 1 |
| 23 | 1 | rw- | 1 | 1 |

A Page Table For 16KB Address Space

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.52 |
|---|---|---|

## PAGE TABLES: WASTED SPACE

- **Process: 16KB Address Space w/ 1KB pages**



**Most of the page table is unused and full of wasted space. (73%)**

| PFN | valid | prot | present | dirty |
|-----|-------|------|---------|-------|
|     |       |      |         | 0     |
|     |       |      |         | -     |
|     |       |      |         | -     |
|     |       |      |         | -     |
| 15  | 1     | rw-  | 1       | 1     |
| ... | ...   |      | ...     | ...   |
| -   | 0     | -    | -       | -     |
| 3   | 1     | rw-  | 1       | 1     |
| 23  | 1     | rw-  | 1       | 1     |

**A Page Table For 16KB Address Space**

**A 16KB Address Space with 1KB Pages**

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.53 |
|---|---|---|

## MULTI-LEVEL PAGE TABLES

- **Consider a page table:**
- **32-bit addressing, 4KB pages**
- **$2^{20}$ page table entries**
- **Even if memory is sparsely populated the *per process* page table requires:**

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

- **Often most of the 4MB *per process* page table is empty**
- **Page table must be placed in 4MB contiguous block of RAM**

- **MUST SAVE MEMORY!**

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.54 |
|---|---|---|

# MULTI-LEVEL PAGE TABLES - 2

- Add level of indirection, the "page directory"



Linear (Left) And Multi-Level (Right) Page Tables

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.55 |

# MULTI-LEVEL PAGE TABLES - 2

- Add level of indirection, the "page directory"



**Two level page table:**
$2^{20}$ **pages addressed with**
**two level-indexing**
**(page directory index, page table index)**

Linear (Left) And Multi-Level (Right) Page Tables

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.56 |

## MULTI-LEVEL PAGE TABLES - 3

- **Advantages**
  - Only allocates page table space in proportion to the address space actually used
  - Can easily grab next free page to expand page table

- **Disadvantages**
  - Multi-level page tables are an example of a time-space tradeoff
  - Sacrifice address translation time (now 2-level) for space
  - Complexity: multi-level schemes are more complex

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.57 |
|---|---|---|

## EXAMPLE

- **16KB address space, 64byte pages**
- **How large would a one-level page table need to be?**
- $2^{14}$ **(address space) / $2^6$ (page size) = $2^8$ = 256 (pages)**



| Flag | Detail |
|---|---|
| Address space | 16 KB |
| Page size | 64 byte |
| Virtual address | 14 bit |
| VPN | 8 bit |
| Offset | 6 bit |
| Page table entry | $2^8$(256) |

**A 16-KB Address Space With 64-byte Pages**

| 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Offset

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.58 |
|---|---|---|

# EXAMPLE - 2

- 256 total page table entries (64 bytes each)

- 1,024 bytes page table size, stored using 64-byte pages
  = (1024/64) = 16 page directory entries (PDEs)

- Each page directory entry (PDE) can hold 16 page table entries (PTEs)  *e.g. lookups*

- 16 page directory entries (PDE) x 16 page table entries (PTE)
  = 256 total PTEs

- **Key idea: the page table is stored using pages too!**

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.59 |
| --- | --- | --- |

# PAGE DIRECTORY INDEX

- Now, let's split the page table into two:
  - 8 bit VPN to map 256 pages
  - 4 bits for <u>page directory index</u> (PDI – 1st level page table)
  - 6 bits offset into 64-byte page



| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.60 |
| --- | --- | --- |

# PAGE TABLE INDEX

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



- 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 | TCSS422: Operating Systems [Fall 2018]<br>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) ?**
- <u>Page directory</u> = 16 entries x 4 bytes (1 x 64 byte page)
- <u>Page table</u> = 4 entries x 4 bytes (1 x 64 byte page)
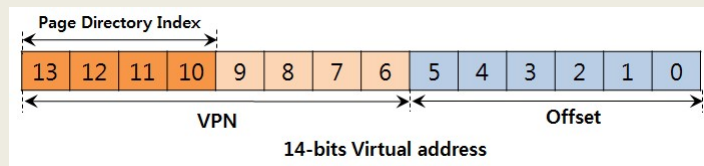- 128 bytes required (2 x 64 byte pages)
  - Savings = using just 12.5% the space !!!

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.62 |
|---|---|---|

## 32-BIT EXAMPLE

- Consider: 32-bit address space, 4KB pages, $2^{20}$ 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** = $2^{10}$ 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

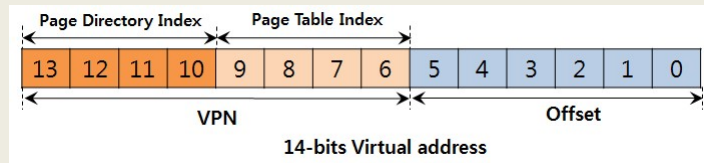| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.63 |
|---|---|---|

## MORE THAN TWO LEVELS

- Consider: page size is $2^9$ = 512 bytes
- Page size 512 bytes / Page entry size 4 bytes
- VPN is 21 bits



| Flag | Detail |
|---|---|
| Virtual address | 30 bit |
| Page size | 512 byte |
| VPN | 21 bit |
| Offset | 9 bit |

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.64 |
|---|---|---|

# MORE THAN TWO LEVELS - 2

- Page table entries per page = 512 / 4 = 128
- 7 bytes – for page table index (PTI)



| Flag | Detail | |
|---|---|---|
| Virtual address | 30 bit | |
| Page size | 512 byte | |
| VPN | 21 bit | |
| Offset | 9 bit | |
| Page entry per page | 128 PTEs | $\rightarrow \log_2 128 = 7$ |

# MORE THAN TWO LEVELS - 3

- To map 1 GB address space ($2^{30}$=1GB RAM, 512-byte pages)
- $2^{14}$ = 16,384 page directory entries (PDEs) are required
- When using $2^7$ (128 entry) page tables…
- Page size = 512 bytes / 4 bytes per addr



| Flag | Detail | |
|---|---|---|
| Virtual address | 30 bit | |
| Page size | 512 byte | |
| VPN | 21 bit | |
| Offset | 9 bit | |
| Page entry per page | 128 PTEs | $\rightarrow \log_2 128 = 7$ |

## MORE THAN TWO LEVELS - 3

- To map 1 GB address space ($2^{30}$=1GB RAM, 512-byte pages)
- $2^{14}$ = 16,384 page directory entries (PDEs) are required
- When using $2^7$ (128 entry) page tables…
- Page size = 512 bytes / 4 bytes per addr

**Can't Store Page Directory with 16K pages, using 512 bytes pages.
Pages only dereference 128 addresses
(512 bytes / 32 bytes)**

| | |
|---|---|
| Virtual address | 30 bit |
| Page size | 512 byte |
| VPN | 21 bit |
| Offset | 9 bit |
| Page entry per page | 128 PTEs |

$\log_2 128 = 7$

November 26, 2018 | TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma | L16.67
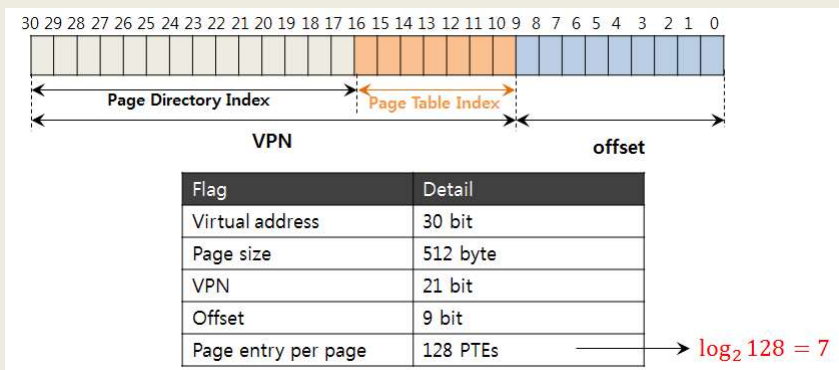
## MORE THAN TWO LEVELS - 3

- To map 1 GB address space ($2^{30}$=1GB RAM, 512-byte pages)
- $2^{14}$ = 16,384 page directory entries (PDEs) are required
- When using $2^7$ (128 entry) page tables…
- Page size = 512 bytes / 4 bytes per addr

**Need three level page table:
Page directory 0 (PD Index 0)
Page directory 1 (PD Index 1)
Page Table Index**

| | |
|---|---|
| Virtual address | 30 bit |
| Page size | 512 byte |
| VPN | 21 bit |
| Offset | 9 bit |
| Page entry per page | 128 PTEs |

$\log_2 128 = 7$

November 26, 2018 | TCSS422: Operating Systems [Fall 2018]
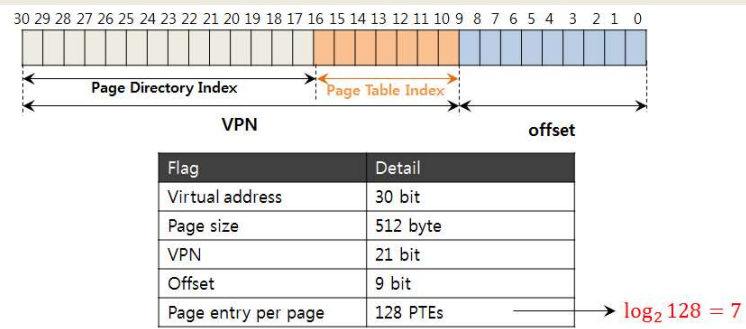School of Engineering and Technology, University of Washington - Tacoma | L16.68

## MORE THAN TWO LEVELS - 4

- We can now address 1GB with "fine grained" 512 byte pages
- Using multiple levels of indirection



- Consider the implications for address translation!
- How much space is required for a virtual address space with 4 entries on a 512-byte page? (let's say 4 32-bit integers)
- PD0 1 page, PD1 1 page, PT 1 page = 1,536 bytes
- Savings = 1,536 / 8,388,608 (8mb) = .0183% !!!

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.69 |

## 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 *pmt;
pte_t *pte;
struct page *page;
```

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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
```

**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

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.71 |
|---|---|---|

## 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 $2^{20}$ pages
- Hash table: can index memory and speed lookups

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.73 |
|---|---|---|

# 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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.74 |
|---|---|---|

# 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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.75 |
|---|---|---|

# 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 | TCSS422: Operating Systems [Fall 2018]<br>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 x 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 → 3.125%

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.77 |

# CHAPTER 21/22: BEYOND PHYSICAL MEMORY

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: Operating Systems [Fall 2018]<br>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**



**Physical Memory and Swap Space**

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

# REPLACEMENT POLICIES

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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 ns, $T_D$ = 10ms
- 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]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.87 |
|---|---|---|

# OPTIMAL REPLACEMENT POLICY

- What if:
  - We could predict the future (… with a magical oracle)
  - All future page accesses are known
  - Always replace the page in the cache used farthest in the future

- Used for a comparison
- Provides a "best case" replacement policy

- Consider a 3-element empty cache with the following page accesses:

0  1  2  0  1  3  0  3  1  2  1

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

**6 hits**

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.88 |
|---|---|---|

# FIFO REPLACEMENT

- Queue based
- Always replace the oldest element at the back of cache
- Simple to implement
- Doesn't consider importance… just arrival ordering

- Consider a 3-element empty cache with the following page accesses:

  0  1  2  0  1  3  0  3  1  2  1

- What is the hit/miss ratio?  **4 hits**
- How is FIFO different than LRU?  **LRU incorporates history**

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.89 |
|---|---|---|

# RANDOM REPLACEMENT

- Pick a page at random to replace
- Simple and fast implementation
- Performance depends on luck of random choices

  0  1  2  0  1  3  0  3  1  2  1



**Random Performance over 10,000 Trials**

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.90 |
|---|---|---|

# 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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.91 |

# WORKLOAD EXAMPLES: NO-LOCALITY

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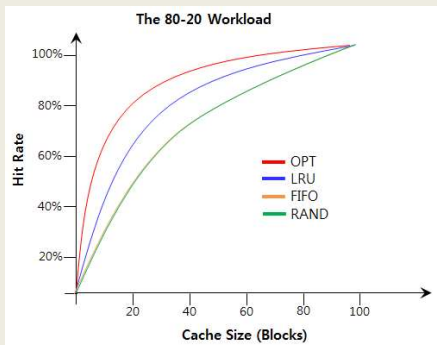  - Perform 10,000 random page accesses
  - Across set of 100 memory pages



**When the cache is large enough to fit the entire workload, it doesn't matter which policy you use.**
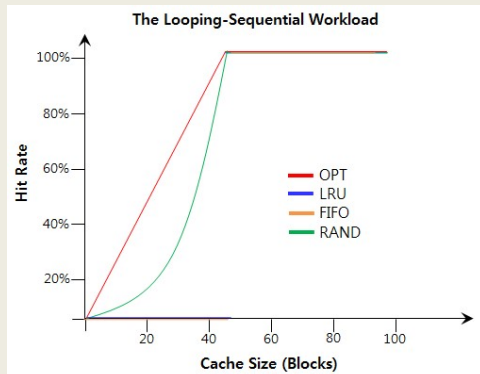
| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.92 |

# 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 | TCSS422: Operating Systems [Fall 2018]<br>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]<br>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]<br>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]<br>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



The 80-20 Workload

| November 26, 2018 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: Operating Systems [Fall 2018]<br>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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.100 |
|---|---|---|

## 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 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L16.101 |
|---|---|---|

# QUESTIONS