


TCCS 422: OPERATING SYSTEMS

Smaller Page Tables,
Beyond Physical Memory



Wes J. Lloyd

School of Engineering and Technology
University of Washington - Tacoma

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

Tacoma

OBJECTIVES

- Mon 3/11 (4pm): Husky Alumni Visit from T-Mobile Q&A
CS work life after graduation–room 206C
Garrett Lahmann ('18), Vlad Kaganyuk ('17)
- Wed 3/13: Prof. Mohamed Ali- UWT CSS Grad Program
- Active Reading Quiz Posted – Chapter 19
- Assignment 3
- Memory Virtualization
 - Chapter 20 – Smaller Page Tables
 - Chapter 21/22 – Beyond Physical Memory

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.2

FEEDBACK FROM 3/6

- Can pages in different locations be used?
- For example: a 48-bytes program, with 3 x 16-byte free pages in different locations?
- Question may refer to example in following slide...

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.3

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

0
16
32
48
64
128

page 0 of the address space (page 1)
page 2
page 3

0
16
32
48
64
80
96
112
128

reserved for OS
(unused)
page 3 of AS
page 0 of AS
(unused)
page 2 of AS
(unused)
page 1 of AS

A Simple 64-byte Address Space

64-Byte Address Space Placed In Physical Memory

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.4

FEEDBACK - 2

- Could you go over the TLB example again?
- Why are they (the array accesses) hits or misses?

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.5

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

00 04 08 12 16

VPN - 00
VPN - 01
VPN - 02
VPN - 03
VPN - 04
VPN - 05
VPN - 06
VPN - 07
VPN - 08
VPN - 09
VPN - 10
VPN - 11
VPN - 12
VPN - 13
VPN - 14
VPN - 15

a[0] a[1] a[2]
a[3] a[4] a[5] a[6]
a[7] a[8] a[9]

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.6

Slides by Wes J. Lloyd

L15.1

TLB EXAMPLE - 2

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

- 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?

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.7

TLB EXAMPLE - 3

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

- 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)

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.8

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

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.9

FEEDBACK - 3


- Still unsure about level 1 & 2 pg table calculations
- Can we do more examples in class?

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.10

CHAPTER 20:
PAGING:
SMALLER TABLES



March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.11

OBJECTIVES

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.12

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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.13

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

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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.14

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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.15

PAGING: USE LARGER PAGES

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

$\frac{2^{32}}{2^{14}} * 4 = 1MB$ 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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.16

PAGE TABLES: WASTED SPACE

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

Page Table

Virtual Address Space

Physical Memory

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

A 16KB Address Space with 1KB Pages

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.17

PAGE TABLES: WASTED SPACE

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

Page Table

Virtual Address Space

Physical Memory

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

PFN	valid	prot	present	dirty
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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.18

MULTI-LEVEL PAGE TABLES

- Consider a page table:
- 32-bit addressing, 4KB pages
- 2²⁰ page table entries
- Even if memory is sparsely populated the *per process* page table requires:

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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.19

MULTI-LEVEL PAGE TABLES - 2

- Add level of indirection, the “page directory”

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.20

MULTI-LEVEL PAGE TABLES - 2

- Add level of indirection, the “page directory”

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.21

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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.22

EXAMPLE

- 16KB address space, 64byte pages
- How large would a one-level page table need to be?
- 2¹⁴ (address space) / 2⁶ (page size) = 2⁸ = 256 (pages)

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.23

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!

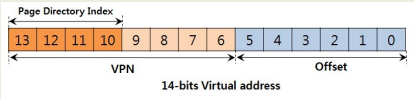
March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.24

PAGE DIRECTORY INDEX

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



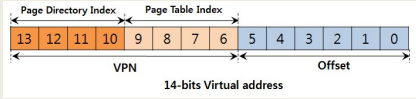
March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.25

PAGE TABLE INDEX

- 4 bits page directory index (PDI – 1st level)
- 4 bits page table index (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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.26

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 !!!

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.27

32-BIT EXAMPLE

- Consider: 32-bit address space, 4KB pages, 2²⁰ 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 = 2¹⁰ 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

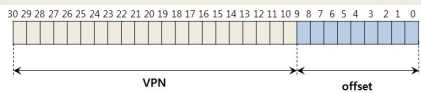
March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.28

MORE THAN TWO LEVELS

- Consider: page size is 2⁹ = 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

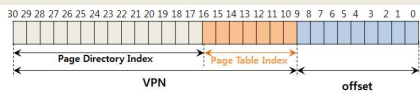
March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.29

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

→ log₂ 128 = 7

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.30

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

→ $\log_2 128 = 7$

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.31

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

→ $\log_2 128 = 7$

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.32

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

→ $\log_2 128 = 7$

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.33

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
- Memory Usage= $1,536 (3\text{-level}) / 8,388,608 (1\text{-level}) = .0183\% !!!$

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.34

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;
```

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.35

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.36

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.37

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?

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.38

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)?

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.39

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.40

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.41

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 \times 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- $512/16384 = .03125 \rightarrow 3.125\%$

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.42

CHAPTER 21/22:
BEYOND PHYSICAL
MEMORY



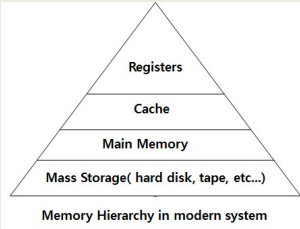
March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.43

MEMORY HIERARCHY

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



Registers
Cache
Main Memory
Mass Storage(hard disk, tape, etc...)
Memory Hierarchy in modern system

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.44

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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.45

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>

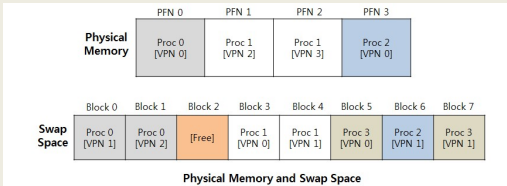
March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.46

SWAP SPACE

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



Physical Memory and Swap Space

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.47

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

March 11, 2019

TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.48

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()
```

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.49

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.50

REPLACEMENT POLICIES

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.51

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 $P_{hit} = .9\text{ (90\%)}$, $P_{miss} = .1$
- Consider $P_{hit} = .999\text{ (99.9\%)}$, $P_{miss} = .001$

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.52

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.53

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.54

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.55

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.56

WORKLOAD EXAMPLES: NO-LOCALITY

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

The No-Locality Workload

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.57

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)

The 80-20 Workload

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.58

WORKLOAD EXAMPLES: SEQUENTIAL

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

The Looping-Sequential Workload

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

Algorithms should provide "scan resistance"

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.59

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.60

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



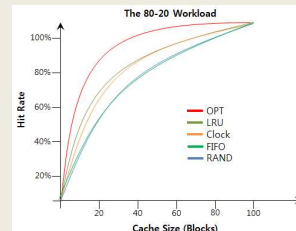
March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.61

CLOCK ALGORITHM

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



March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.62

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.63

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.64

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.65

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

March 11, 2019

TCCS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L15.66

