


TCCS 422: OPERATING SYSTEMS

Smaller Tables,  
Beyond Virtual Memory



Wes J. Lloyd

Institute of Technology

University of Washington - Tacoma

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

OBJECTIVES

- Homework 3 Questions
- Ch. 20
  - Smaller Tables
- Ch. 21/22
  - Beyond Physical Memory: Mechanisms (Ch. 21)
    - Virtual "Swap" Memory
  - Beyond Physical Memory: Policies (Ch. 22)
    - Page Replacement Algorithms
    - Replacement algorithm effectiveness
- Ch. 36
  - I/O Devices

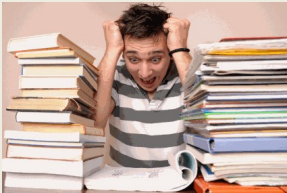
March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.2

FEEDBACK FROM 2/28

- There was no feedback !!!




March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.3

CHAPTER 20:  
PAGING:  
SMALLER TABLES



March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.4

MULTI-LEVEL PAGE TABLES

- Consider a page table:
  - 32-bit addressing, 4KB pages
  - 2<sup>20</sup> page table entries
- Even if memory is sparsely populated the *per process* page table requires:

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!

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.5

MULTI-LEVEL PAGE TABLES - 2

- Add level of indirection, the "page directory"

Linear Page Table

valid	prot	PFN
1	rw	12
1	rw	13
0	-	-
1	rw	100

Multi-level Page Table

valid	PFN
1	201
0	-
1	203

The Page Directory

valid	prot	PFN
1	rw	12
1	rw	13
0	-	-
1	rw	100

[Page 1 of PT Not Allocated]

valid	prot	PFN
0	-	-
0	-	-
1	rw	86
1	rw	15

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

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.6

MULTI-LEVEL PAGE TABLES - 2

- Add level of indirection, the "page directory"

Linear Page Table

Multi-level Page Table

Two level page table:  
2<sup>20</sup> pages addressed with  
two level-indexing  
(page directory index, page table index)

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

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.7

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.8

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

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.9

MORE THAN TWO LEVELS

- Consider: page size is 2<sup>9</sup> = 512 bytes
- Page size 512 bytes / Page entry size 4 bytes
- VPN is 21 bits

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

VPN

offset

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

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.10

MORE THAN TWO LEVELS - 2

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

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

Page Directory Index

Page Table Index

offset

VPN

offset

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

→ log<sub>2</sub> 128 = 7

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.11

MORE THAN TWO LEVELS - 3

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

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

Page Directory Index

Page Table Index

offset

VPN

offset

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

→ log<sub>2</sub> 128 = 7

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.12

Slides by Wes J. Lloyd

L15.2

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 address

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.13

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 address

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.14

MORE THAN TWO LEVELS - 4

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

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

PD Index 0 PD Index 1 Page Table Index

VPN

- 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% !!!

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.15

ADDRESS TRANSLATION - 1

```
01: VPN = (VirtualAddress & VPN_MASK) >> SHIFT
02: (Success, TlbEntry) = TLB_Lookup(VPN)
03: if (Success == True) //TLB Hit
04:     if (CanAccess(TlbEntry.ProtectBits) == True)
05:         Offset = VirtualAddress & OFFSET_MASK
06:         PhysAddr = (TlbEntry.PFN << SHIFT) | Offset
07:         Register = AccessMemory(PhysAddr)
08:     else RaiseException(PROTECTION_FAULT);
09: else // perform the full multi-level lookup
```

(05-07) Generate physical address from TLB

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.16

ADDRESS TRANSLATION - 2

```
11: else
12:     PDIndex = (VPN & PD_MASK) >> PD_SHIFT
13:     PDEAddr = PDBase + (PDIndex * sizeof(PDE))
14:     PDE = AccessMemory(PDEAddr)
15:     if (PDE.Valid == False)
16:         RaiseException(SEGMENTATION_FAULT)
17:     else // PDE is Valid: now fetch PTE from PT
```

(15-17) Check if PDE is valid, if so fetch entry from page table

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.17

ADDRESS TRANSLATION - 3


```
18: PTIndex = (VPN & PT_MASK) >> PT_SHIFT
19: PTEAddr = (PDE.PFN << SHIFT) + (PTIndex * sizeof(PTE))
20: PTE = AccessMemory(PTEAddr)
21: if (PTE.Valid == False)
22:     RaiseException(SEGMENTATION_FAULT)
23: else if (CanAccess(PTE.ProtectBits) == False)
24:     RaiseException(PROTECTION_FAULT);
25: else
26:     TLB_Insert(VPN, PTE.PFN, PTE.ProtectBits)
27:     RetryInstruction()
```

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.18

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.19

CHAPTER 21/22:  
BEYOND PHYSICAL  
MEMORY



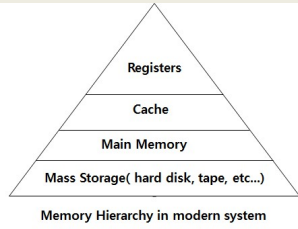
March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.20

MEMORY HIERARCHY

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



March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.21

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.22

LATENCY TIMES

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

Action	Latency (ns)	( $\mu$ s)	
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 $\mu$ s	~1GB/sec SSD
Read 1 MB sequentially from memory	250,000 ns	250 $\mu$ s	
Read 1 MB sequentially from SSD*	1,000,000 ns	1,000 $\mu$ s	1 ms ~1GB/sec SSD, 4X memory
Read 1 MB sequentially from disk	20,000,000 ns	20,000 $\mu$ s	20 ms 80x memory, 20X SSD

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

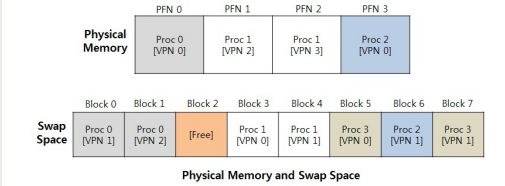
March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.23

SWAP SPACE

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



March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.24

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.25

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
```

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.26

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.27

REPLACEMENT POLICIES



March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.28

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%),  $P_{miss} = .001$

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.29

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.30

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 5, 2018

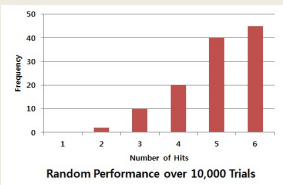
TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.31

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 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.32

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

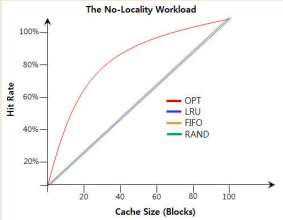
March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.33

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.

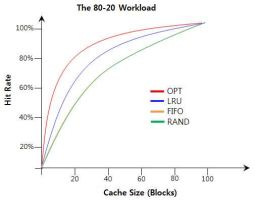
March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.34

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)

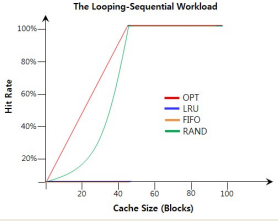
March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.35

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"

March 5, 2018

TCCS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.36

## 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 5, 2018

TCS5422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.37

## 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



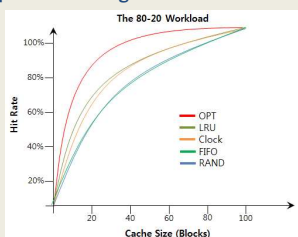
March 5, 2018

TCS5422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.38

## CLOCK ALGORITHM

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



March 5, 2018

TCS5422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.39

## 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 5, 2018

TCS5422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.40

## 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 5, 2018

TCS5422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.41

## 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 5, 2018

TCS5422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.42

## 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 5, 2018

TCSS422: Operating Systems [Winter 2018]  
Institute of Technology, University of Washington - Tacoma

L15.43

## QUESTIONS

