# TCSS 422: OPERATING SYSTEMS

## Beyond Physical Memory

Wes J. Lloyd
Institute of Technology
University of Washington - Tacoma

---

## OBJECTIVES

- Chapter 21
  - Virtual "Swap" Memory

- Chapter 22
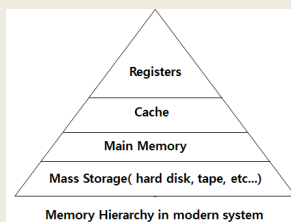  - Page replacement algorithms

  - Replacement algorithm effectiveness

| March 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.2 |

---

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.3 |

---

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.4 |

---

## 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 |

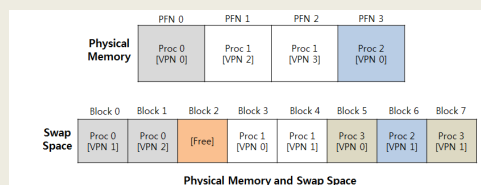- Latency numbers every programmer should know
- From: https://gist.github.com/jboner/2841832#file-latency-txt

| March 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.5 |

---

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.6 |

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.7

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.8

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.9

## REPLACEMENT POLICIES

March 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.10

## 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

March 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.11

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.12

## FIFO REPLACEMENT

- Queue based
- Always replace the oldest element
- 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**
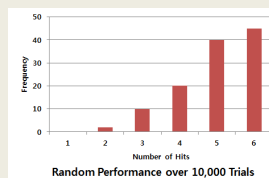- How is FIFO different than LRU?   **LRU incorporates history**

| March 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.13 |

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.14 |

## HISTORY-BASED POLICIES

- LRU: Least recently used
- Always replace page with oldest access time
- Consider when a page 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
- Consider frequency of page accesses
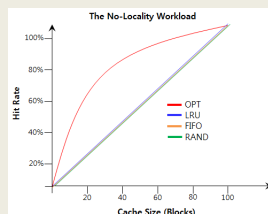
  0 1 2 0 1 3 0 3 1 2 1   **Hit/miss ratio is=**
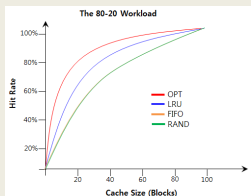                          **6 hits**

| March 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.15 |

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.16 |

## 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

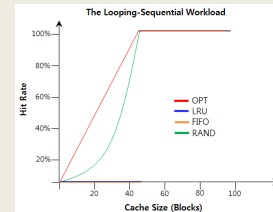LRU is more likely to hold onto hot pages

(recalls history)

| March 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.17 |

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017]<br>Institute of Technology, University of Washington - Tacoma | L17.18 |

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.19 |

## 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
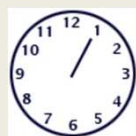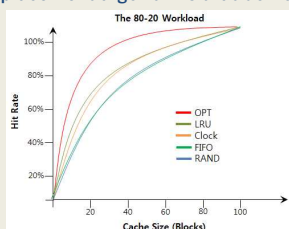    - IF USE_BIT=0 replace page

| March 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.20 |

## CLOCK ALGORITHM

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



| March 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.21 |

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.22 |

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.23 |

## 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 6, 2017 | TCSS422: Operating Systems [Winter 2017] Institute of Technology, University of Washington - Tacoma | L17.24 |

## 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 6, 2017    TCSS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma    L17.25

# QUESTIONS

March 6, 2017    TCSS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma    L17.26