# CS162
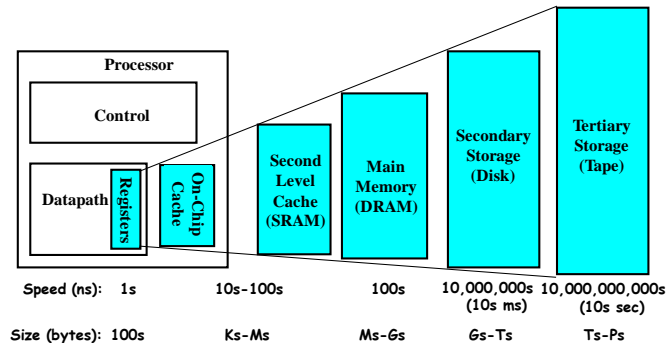# Operating Systems and Systems Programming
# Lecture 14

# Caching and Demand Paging

March 17, 2008
Prof. Anthony D. Joseph
http://inst.eecs.berkeley.edu/~cs162

---

## Review: Memory Hierarchy of a Modern Computer System

- Take advantage of the principle of locality to:
  - Present as much memory as in the cheapest technology
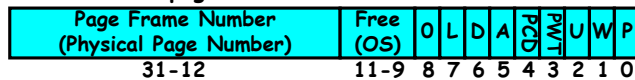  - Provide access at speed offered by the fastest technology



| | Processor | Second Level Cache (SRAM) | Main Memory (DRAM) | Secondary Storage (Disk) | Tertiary Storage (Tape) |
|---|---|---|---|---|---|
| Speed (ns): | 1s | 10s-100s | 100s | 10,000,000s (10s ms) | 10,000,000,000s (10s sec) |
| Size (bytes): | 100s | Ks-Ms | Ms-Gs | Gs-Ts | Ts-Ps |

---

## Review: What is in a PTE?

- What is in a Page Table Entry (or PTE)?
  - Pointer to next-level page table or to actual page
  - Permission bits: valid, read-only, read-write, write-only
- Example: Intel x86 architecture PTE:
  - Address same format previous slide (10, 10, 12-bit offset)
  - Intermediate page tables called "Directories"

| Page Frame Number (Physical Page Number) | Free (OS) | 0 | L | D | A | PCD | PWT | U | W | P |
|---|---|---|---|---|---|---|---|---|---|---|
| 31-12 | 11-9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

  P: Present (same as "valid" bit in other architectures)
  W: Writeable
  U: User accessible
PWT: Page write transparent: external cache write-through
PCD: Page cache disabled (page cannot be cached)
  A: Accessed: page has been accessed recently
  D: Dirty (PTE only): page has been modified recently
  L: L=1⇒4MB page (directory only).
     Bottom 22 bits of virtual address serve as offset

---

## Review: Other Caching Questions

- What line gets replaced on cache miss?
  - Easy for Direct Mapped: Only one possibility
  - Set Associative or Fully Associative:
    » Random
    » LRU (Least Recently Used)
- What happens on a write?
  - Write through: The information is written to both the cache and to the block in the lower-level memory
  - Write back: The information is written only to the block in the cache
    » Modified cache block is written to main memory only when it is replaced
    » Question is block clean or dirty?

Page 1

## Goals for Today

- Concept of Paging to Disk
- Page Faults and TLB Faults
- Precise Interrupts
- Page Replacement Policies

Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne. Many slides generated from my lecture notes by Kubiatowicz.
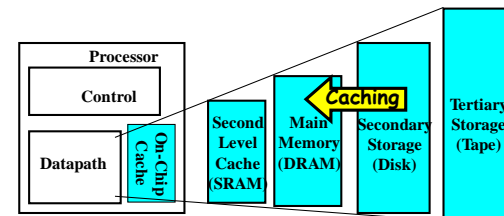
## Demand Paging

- Modern programs require a lot of physical memory
  - Memory per system growing faster than 25%-30%/year
- But they don't use all their memory all of the time
  - 90-10 rule: programs spend 90% of their time in 10% of their code
  - Wasteful to require all of user's code to be in memory
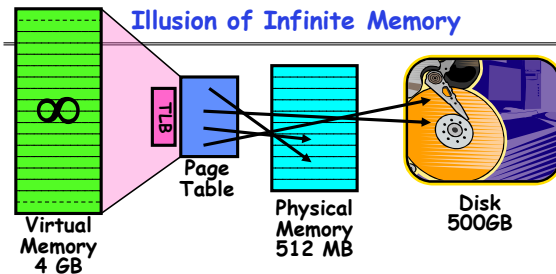- Solution: use main memory as cache for disk

## Illusion of Infinite Memory



- Disk is larger than physical memory ⇒
  - In-use virtual memory can be bigger than physical memory
  - Combined memory of running processes much larger than physical memory
    - » More programs fit into memory, allowing more concurrency
- Principle: Transparent Level of Indirection (page table)
  - Supports flexible placement of physical data
    - » Data could be on disk or somewhere across network
  - Variable location of data transparent to user program
    - » Performance issue, not correctness issue

## Demand Paging is Caching

- Since Demand Paging is Caching, must ask:
  - What is block size?
    - » 1 page
  - What is organization of this cache (i.e. direct-mapped, set-associative, fully-associative)?
    - » Fully associative: arbitrary virtual→physical mapping
  - How do we find a page in the cache when look for it?
    - » First check TLB, then page-table traversal
  - What is page replacement policy? (i.e. LRU, Random…)
    - » This requires more explanation… (kinda LRU)
  - What happens on a miss?
    - » Go to lower level to fill miss (i.e. disk)
  - What happens on a write? (write-through, write back)
    - » Definitely write-back. Need dirty bit!

## Administrative

- **Project 2 code due Thursday 3/20 at 11:59pm**
  - Project 2 autograder is up and running every 15 minutes

- **Make sure you attend sections!**
  - There will be a lot of information about the projects that I cannot cover in class
  - Also supplemental information and detail that we don't have time for in class

- **We have an anonymous feedback link on the course homepage**
  - Please use to give feedback on course
  - Wednesday: We will have a survey to fill out

## Demand Paging Mechanisms

- **PTE helps us implement demand paging**
  - Valid ⇒ Page in memory, PTE points at physical page
  - Not Valid ⇒ Page not in memory; use info in PTE to find it on disk when necessary
- **Suppose user references page with invalid PTE?**
  - Memory Management Unit (MMU) traps to OS
    - » Resulting trap is a "Page Fault"
  - What does OS do on a Page Fault?:
    - » Choose an old page to replace
    - » If old page modified ("D=1"), write contents back to disk
    - » Change its PTE and any cached TLB to be invalid
    - » Load new page into memory from disk
    - » Update page table entry, invalidate TLB for new entry
    - » Continue thread from original faulting location
  - TLB for new page will be loaded when thread continued!
  - While pulling pages off disk for one process, OS runs another process from ready queue
    - » Suspended process sits on wait queue

## Software-Loaded TLB
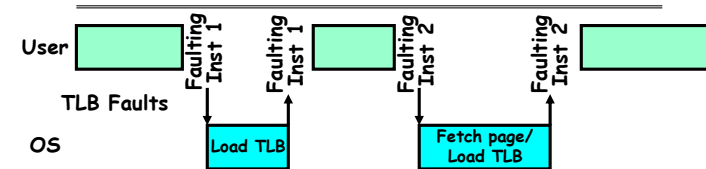
- **MIPS/Nachos TLB is loaded by software**
  - High TLB hit rate⇒ok to trap to software to fill the TLB, even if slower
  - Simpler hardware and added flexibility: software can maintain translation tables in whatever convenient format
- **How can a process run without access to page table?**
  - Fast path (TLB hit with valid=1):
    - » Translation to physical page done by hardware
  - Slow path (TLB hit with valid=0 or TLB miss)
    - » Hardware receives a "TLB Fault"
  - What does OS do on a TLB Fault?
    - » Traverse page table to find appropriate PTE
    - » If valid=1, load page table entry into TLB, continue thread
    - » If valid=0, perform "Page Fault" detailed previously
    - » Continue thread
- **Everything is transparent to the user process:**
  - It doesn't know about paging to/from disk
  - It doesn't even know about software TLB handling

## Transparent Exceptions



- **How to transparently restart faulting instructions?**
  - Could we just skip it?
    - » No: need to perform load or store after reconnecting physical page
- **Hardware must help out by saving:**
  - Faulting instruction and partial state
    - » Need to know which instruction caused fault
    - » Is single PC sufficient to identify faulting position????
  - Processor State: sufficient to restart user thread
    - » Save/restore registers, stack, etc
- **What if an instruction has side-effects?**

## Consider weird things that can happen

- **What if an instruction has side effects?**
  - Options:
    - » Unwind side-effects (easy to restart)
    - » Finish off side-effects (messy!)
  - Example 1: `mov (sp)+,10`
    - » What if page fault occurs when write to stack pointer?
    - » Did `sp` get incremented before or after the page fault?
  - Example 2: `strcpy (r1), (r2)`
    - » Source and destination overlap: can't unwind in principle!
    - » IBM S/370 and VAX solution: execute twice – once read-only
- **What about "RISC" processors?**
  - For instance delayed branches?
    - » Example:    `bne somewhere`
                    `ld r1,(sp)`
    - » Precise exception state consists of two PCs: PC and nPC
  - Delayed exceptions:
    - » Example:    `div r1, r2, r3`
                    `ld r1, (sp)`
    - » What if takes many cycles to discover divide by zero, but load has already caused page fault?
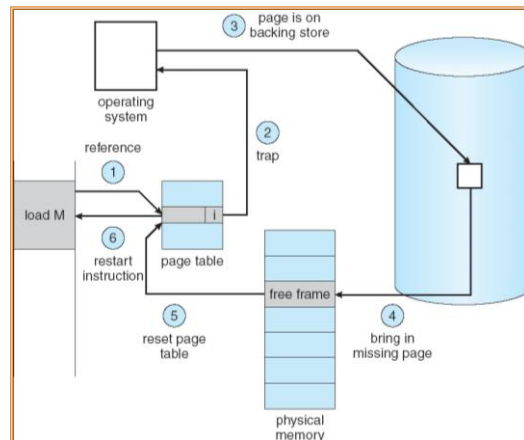
## Precise Exceptions

- **Precise** $\Rightarrow$ **state of the machine is preserved as if program executed up to the offending instruction**
  - All previous instructions **completed**
  - Offending instruction and all following instructions act **as if they have not even started**
  - Same system code will work on different implementations
  - Difficult in the presence of pipelining, out-of-order execution, ...
  - **MIPS takes this position**
- **Imprecise** $\Rightarrow$ **system software has to figure out what is where and put it all back together**
- **Performance goals often lead designers to forsake precise interrupts**
  - system software developers, user, markets etc. usually wish they had not done this
- **Modern techniques for out-of-order execution and branch prediction help implement precise interrupts**

## Steps in Handling a Page Fault

## Demand Paging Example

- **Since Demand Paging like caching, can compute average access time! ("Effective Access Time")**
  - EAT = Hit Rate x Hit Time + Miss Rate x Miss Time
- **Example:**
  - Memory access time = 200 nanoseconds
  - Average page-fault service time = 8 milliseconds
  - Suppose p = Probability of miss, 1-p = Probably of hit
  - Then, we can compute EAT as follows:
    
    EAT  = (1 – p) x 200ns + p x 8 ms
          = (1 – p)  x 200ns + p x 8,000,000ns
          = 200ns + p x 7,999,800ns
- **If one access out of 1,000 causes a page fault, then EAT = 8.2 μs:**
  - This is a slowdown by a factor of 40!
- **What if want slowdown by less than 10%?**
  - 200ns x 1.1 < EAT $\Rightarrow$ p < 2.5 x $10^{-6}$
  - This is about 1 page fault in 400000!

## Review: What Factors Lead to Misses?

- **Compulsory Misses:**
  - Pages that have never been paged into memory before
  - How might we remove these misses?
    - » Prefetching: loading them into memory before needed
    - » Need to predict future somehow! More later…
- **Capacity Misses:**
  - Not enough memory – must somehow increase size
    - » One option: Increase amount of DRAM (not quick fix!)
    - » Another option: If multiple processes in memory: adjust percentage of memory allocated to each one!
- **Conflict Misses (collision):**
  - Doesn't exist in virtual memory ("fully-associative" cache)
- **Policy Misses:**
  - Replacement policy kicks pages out of memory early
  - How to fix? Better replacement policy
- **Coherence Misses (invalidation):**
  - Not a problem for virtual memory, since all processes/cores share same memory

---

**BREAK**

---

## Page Replacement Policies

- **Why do we care about Replacement Policy?**
  - Replacement is an issue with any cache
  - Particularly important with pages
    - » The cost of being wrong is high: must go to disk
    - » Must keep important pages in memory, not toss them out
- **FIFO (First In, First Out)**
  - Throw out oldest page. Be fair – let every page live in memory for same amount of time.
  - Bad, because throws out heavily used pages instead of infrequently used pages
- **MIN (Minimum):**
  - Replace page that won't be used for the longest time
  - Great, but can't really know future…
  - Makes good comparison case, however
- **RANDOM:**
  - Pick random page for every replacement
  - Typical solution for TLB's. Simple hardware
  - Pretty unpredictable – makes it hard to make real-time guarantees

---

## Replacement Policies (Con't)

- **LRU (Least Recently Used):**
  - Replace page that hasn't been used for the longest time
  - Programs have locality, so if something not used for a while, unlikely to be used in the near future.
  - Seems like LRU should be a good approximation to MIN.
- **How to implement LRU? Use a list!**

Head → Page 6 → Page 7 → Page 1 → Page 2

Tail (LRU)

  - On each use, remove page from list and place at head
  - LRU page is at tail
- **Problems with this scheme for paging?**
  - Need to know immediately when each page used so that can change position in list…
  - Many instructions for each hardware access
- **In practice, people approximate LRU (more later)**

## Example: FIFO

- **Suppose we have 3 page frames, 4 virtual pages, and following reference stream:**
  - A B C A B D A D B C B
- **Consider FIFO Page replacement:**

| Ref: | A | B | C | A | B | D | A | D | B | C | B |
|------|---|---|---|---|---|---|---|---|---|---|---|
| **Page:** | | | | | | | | | | | |
| **1** | A | | | | | D | | | | C | |
| **2** | | B | | | | | A | | | | |
| **3** | | | C | | | | | | B | | |

- FIFO: 7 faults.
- When referencing D, replacing A is bad choice, since need A again right away

## Example: MIN

- **Suppose we have the same reference stream:**
  - A B C A B D A D B C B
- **Consider MIN Page replacement:**

| Ref: | A | B | C | A | B | D | A | D | B | C | B |
|------|---|---|---|---|---|---|---|---|---|---|---|
| **Page:** | | | | | | | | | | | |
| **1** | A | | | | | | | | | C | |
| **2** | | B | | | | | | | | | |
| **3** | | | C | | | D | | | | | |

- MIN: 5 faults
- Where will D be brought in? Look for page not referenced farthest in future.
- **What will LRU do?**
  - Same decisions as MIN here, but won't always be true!

## When will LRU perform badly?

- **Consider the following: A B C D A B C D A B C D**
- **LRU Performs as follows (same as FIFO here):**

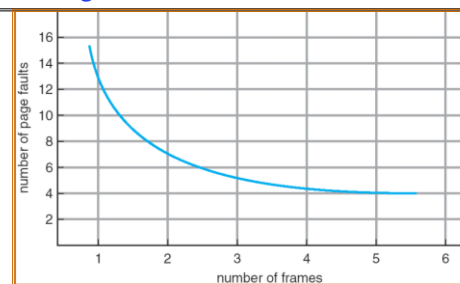| Ref: | A | B | C | D | A | B | C | D | A | B | C | D |
|------|---|---|---|---|---|---|---|---|---|---|---|---|
| **Page:** | | | | | | | | | | | | |
| **1** | A | | | D | | | C | | | B | | |
| **2** | | B | | | A | | | D | | | C | |
| **3** | | | C | | | B | | | A | | | D |

  - Every reference is a page fault!
- **MIN Does much better:**

| Ref: | A | B | C | D | A | B | C | D | A | B | C | D |
|------|---|---|---|---|---|---|---|---|---|---|---|---|
| **Page:** | | | | | | | | | | | | |
| **1** | A | | | | | | | | | B | | |
| **2** | | B | | | | | C | | | | | |
| **3** | | | C | D | | | | | | | | |

## Graph of Page Faults Versus The Number of Frames



- **One desirable property: When you add memory the miss rate goes down**
  - Does this always happen?
  - Seems like it should, right?
- **No: BeLady's anomaly**
  - Certain replacement algorithms (FIFO) don't have this obvious property!

## Adding Memory Doesn't Always Help Fault Rate

- Does adding memory reduce number of page faults?
  - Yes for LRU and MIN
  - Not necessarily for FIFO!  (Called Belady's anomaly)

| Ref: Page: | A | B | C | D | A | B | E | A | B | C | D | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | A |   |   | D |   |   | E |   |   |   |   |   |
| 2 |   | B |   |   | A |   |   |   |   | C |   |   |
| 3 |   |   | C |   |   | B |   |   |   |   | D |   |

| Ref: Page: | A | B | C | D | A | B | E | A | B | C | D | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | A |   |   |   |   |   | E |   |   |   | D |   |
| 2 |   | B |   |   |   |   |   | A |   |   |   | E |
| 3 |   |   | C |   |   |   |   |   | B |   |   |   |
| 4 |   |   |   | D |   |   |   |   |   | C |   |   |

- After adding memory:
  - With FIFO, contents can be completely different
  - In contrast, with LRU or MIN, contents of memory with X pages are a subset of contents with X+1 Page

## Implementing LRU

- Perfect:
  - Timestamp page on each reference
  - Keep list of pages ordered by time of reference
  - Too expensive to implement in reality for many reasons
- Clock Algorithm: Arrange physical pages in circle with single clock hand
  - Approximate LRU (approx to approx to MIN)
  - Replace an old page, not the oldest page
- Details:
  - Hardware "use" bit per physical page:
    » Hardware sets use bit on each reference
    » If use bit isn't set, means not referenced in a long time
    » Nachos hardware sets use bit in the TLB; you have to copy this back to page table when TLB entry gets replaced
  - On page fault:
    » Advance clock hand (not real time)
    » Check use bit: 1→used recently; clear and leave alone
       0→selected candidate for replacement
  - Will always find a page or loop forever?
    » Even if all use bits set, will eventually loop around⇒FIFO

Set of all pages in Memory

## Summary

- Demand Paging:
  - Treat memory as cache on disk
  - Cache miss ⇒ get page from disk
- Transparent Level of Indirection
  - User program is unaware of activities of OS behind scenes
  - Data can be moved without affecting application correctness
- Software-loaded TLB
  - Fast Path: handled in hardware (TLB hit with valid=1)
  - Slow Path: Trap to software to scan page table
- Precise Exception specifies a single instruction for which:
  - All previous instructions have completed (committed state)
  - No following instructions nor actual instruction have started
- Replacement policies
  - FIFO: Place pages on queue, replace page at end
  - MIN: replace page that will be used farthest in future
  - LRU: Replace page that hasn't be used for the longest time
  - Clock Algorithm: Approximation to LRU

Page 7