# CS162
# Operating Systems and
# Systems Programming
# Lecture 11

# Page Allocation and Replacement

February 28, 2011
Ion Stoica
http://inst.eecs.berkeley.edu/~cs162

---

## Goals for Today

- Finish discussion on TLBs
- Page Replacement Policies
  - FIFO, LRU
  - Clock Algorithm
- Working Set/Thrashing

**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.**
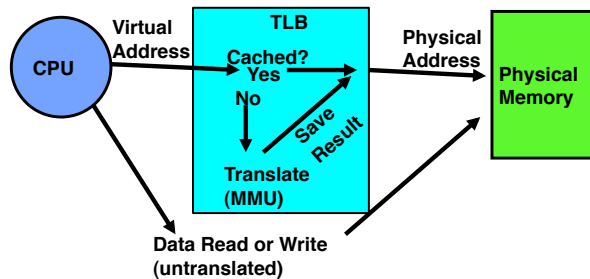
---

## Review: Caching Applied to Address Translation

- Problem: address translation expensive (especially multi-level)
- Solution: cache address translation (TLB)
  - Instruction accesses spend a lot of time on the same page (since accesses sequential)
  - Stack accesses have definite locality of reference
  - Data accesses have less page locality, but still some…



**CPU** → **Virtual Address** → **TLB** Cached? Yes / No → **Physical Address** → **Physical Memory**

Translate (MMU)

Save Result

Data Read or Write (untranslated)

---

## TLB organization

- How big does TLB actually have to be?
  - Usually small: 128-512 entries
  - Not very big, can support higher associativity

- TLB usually organized as fully-associative cache
  - Lookup is by Virtual Address
  - Returns Physical Address

- What happens when fully-associative is too slow?
  - Put a small (4-16 entry) direct-mapped cache in front
  - Called a "TLB Slice"

- When does TLB lookup occur?
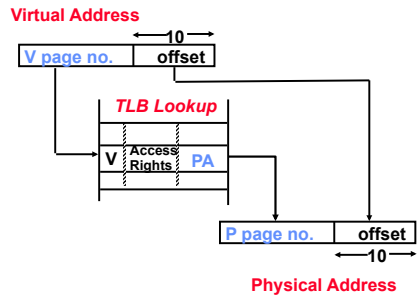  - Before cache lookup?
  - In parallel with cache lookup?

---

Page 1

## Reducing translation time further

- As described, TLB lookup is in serial with cache lookup:

**Virtual Address**

| V page no. | offset |
|---|---|

←10→

**TLB Lookup**

| V | Access Rights | PA |
|---|---|---|

| P page no. | offset |
|---|---|

←10→

**Physical Address**

- Machines with TLBs go one step further: they overlap TLB lookup with cache access.
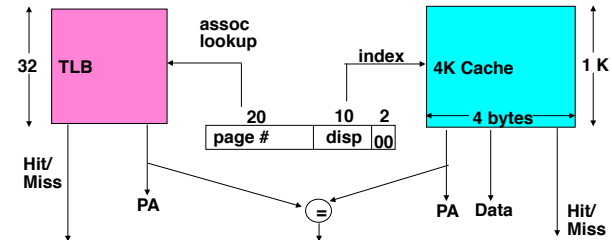  – Works because offset available early

## Overlapping TLB & Cache Access

- Here is how this might work with a 4K cache:



32    TLB    assoc lookup    index    4K Cache    1 K

| 20 | 10 | 2 |
|---|---|---|
| page # | disp | 00 |

4 bytes

Hit/Miss          PA          =          PA    Data          Hit/Miss
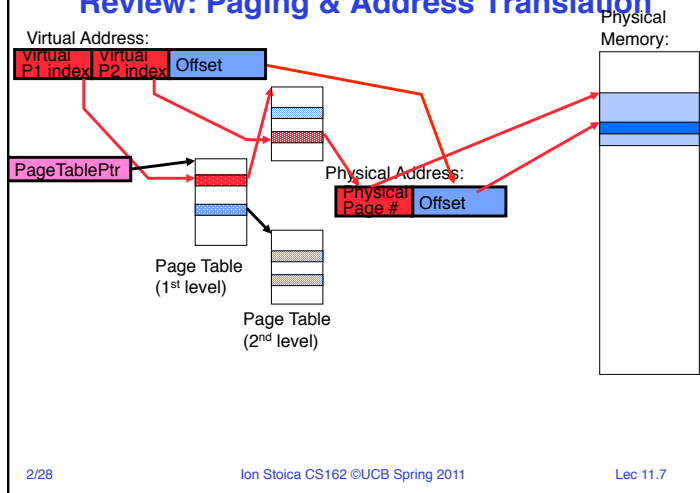
- What if cache size is increased to 8KB?
  – Overlap not complete
  – Need to do something else.  See CS152/252
- Another option: Virtual Caches
  – Tags in cache are virtual addresses
  – Translation only happens on cache misses

## Review: Paging & Address Translation

Virtual Address:

| Virtual P1 index | Virtual P2 index | Offset |
|---|---|---|

PageTablePtr

Physical Memory:

Physical Address:

| Physical Page # | Offset |
|---|---|

Page Table (1st level)

Page Table (2nd level)

## Review: Translation Look-aside Buffer

Virtual Address:

| Virtual P1 index | Virtual P2 index | Offset |
|---|---|---|

PageTablePtr

Physical Memory:

Physical Address:

| Physical Page # | Offset |
|---|---|

Page Table (1st level)

Page Table (2nd level)

TLB:

Page 2

## Review: Cache

Virtual Address:

Virtual P1 index | Virtual P2 index | Offset

PageTablePtr

Physical Address:
Physical Page # | Offset

tag | index | byte

cache:
tag: | block:

Page Table (1st level)

Page Table (2nd level)

TLB:

Physical Memory:

---

## 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

Processor
Control
Datapath
On-Chip Cache
Second Level Cache (SRAM
Main Memory (DRAM)
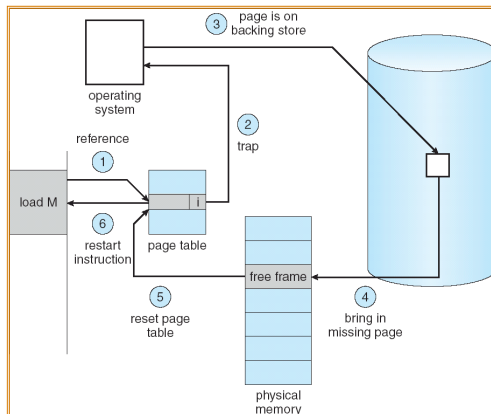Caching
Secondary Storage (Disk)
Tertiary Storage (Tape)

---

## 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 a "dirty" bit (D)!

---

## 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

Page 3

## Steps in Handling a Page Fault



page is on backing store ③

operating system

reference ①   trap ②

load M

restart instruction ⑥   page table   free frame

reset page table ⑤   bring in missing page ④

physical memory

---

## 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
    $\quad\quad = (1 - p)$ x 200ns + p x 8,000,000ns
    $\quad\quad = $ 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 400,000 !

---

## 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.
  - Can we do this?
    » 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:
  - Technically, conflict misses don't exist in virtual memory, since it is a "fully-associative" cache
- Policy Misses:
  - Caused when pages were in memory, but kicked out prematurely because of the replacement policy
  - How to fix? Better replacement policy

---

## 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
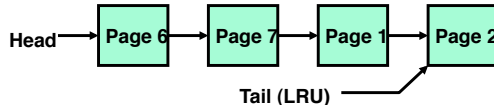  - Typical solution for TLB's. Simple hardware
  - Unpredictable

Page 4

# 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?
  - List operations complex
    » 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
  - 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
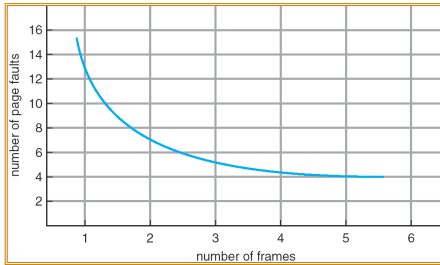- LRU Performs as follows (same as FIFO here):

| 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!

---

## Administrivia

- Project 1
  – Code: Tuesday, March 1st
  – Final document, peer evaluation: Wednesday, March 2nd

- Midterm next week:
  – Wednesday, March 9th
  – Closed book, one page of hand-written notes (both sides)

- Midterm Topics: Everything up to this Wednesday, March 2nd

---

## 5min Break

---

## 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)

| 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 |   |

| 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

Page 6

## Implementing LRU & Second Chance

- Perfect:
  - Timestamp page on each reference
  - Keep list of pages ordered by time of reference
  - Too expensive to implement in reality for many reasons

- Second Chance Algorithm:
  - Approximate LRU
    » Replace an old page, not the oldest page
  - FIFO with "use" bit

- Details
  - A "use" bit per physical page
  - On page fault check page at head of queue
    » If use bit=1 → clear bit, and move page at tail (give the page second chance!)
    » If use bit=0 → replace page
  - Moving pages to tail still complex

## Clock Algorithm

- Clock Algorithm: more efficient implementation of second chance algorithm
  - Arrange physical pages in circle with single clock hand
- Details:
  - On page fault:
    » Advance clock hand (not real time)
    » Check use bit: 1→used recently; clear and leave it alone
                     0→selected candidate for replacement
  - Will always find a page or loop forever?

- What if hand moving slowly?
  - Good sign or bad sign?
    » Not many page faults and/or find page quickly
- What if hand is moving quickly?
  - Lots of page faults and/or lots of reference bits set

## Second Chance Illustration

- Max page table size 4
  - Page B arrives
  - Page A arrives
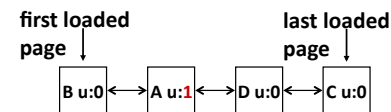  - Access page A
  - Page D arrives
  - Page C arrives

**first loaded page** ↓          **last loaded page** ↓

[ B u:0 ]↔[ A u:**1** ]↔[ D u:0 ]↔[ C u:0 ]

## Second Chance Illustration

- Max page table size 4
  - Page B arrives
  - Page A arrives
  - Access page A
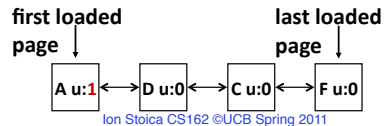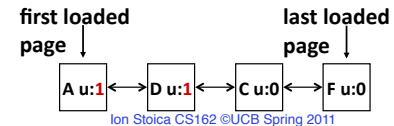  - Page D arrives
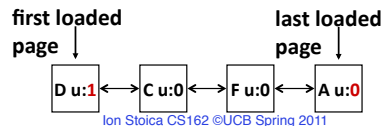  - Page C arrives
  - Page F arrives

**first loaded page** ↓          **last loaded page** ↓

[ B u:0 ]↔[ A u:**1** ]↔[ D u:0 ]↔[ C u:0 ]

## Second Chance Illustration

- Max page table size 4
  - Page B arrives
  - Page A arrives
  - Access page A
  - Page D arrives
  - Page C arrives
  - Page F arrives

**first loaded page** ↓      **last loaded page** ↓

| A u:**1** | ↔ | D u:0 | ↔ | C u:0 | ↔ | F u:0 |

## Second Chance Illustration

- Max page table size 4
  - Page B arrives
  - Page A arrives
  - Access page A
  - Page D arrives
  - Page C arrives
  - Page F arrives
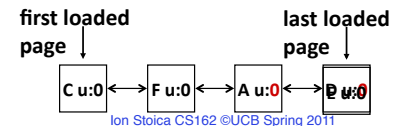  - Access page D
  - Page E arrives

**first loaded page** ↓      **last loaded page** ↓

| A u:**1** | ↔ | D u:**1** | ↔ | C u:0 | ↔ | F u:0 |

## Second Chance Illustration

- Max page table size 4
  - Page B arrives
  - Page A arrives
  - Access page A
  - Page D arrives
  - Page C arrives
  - Page F arrives
  - Access page D
  - Page E arrives

**first loaded page** ↓      **last loaded page** ↓

| D u:**1** | ↔ | C u:0 | ↔ | F u:0 | ↔ | A u:**0** |

## Second Chance Illustration

- Max page table size 4
  - Page B arrives
  - Page A arrives
  - Access page A
  - Page D arrives
  - Page C arrives
  - Page F arrives
  - Access page D
  - Page E arrives

**first loaded page** ↓      **last loaded page** ↓

| C u:0 | ↔ | F u:0 | ↔ | A u:**0** | ↔ | D u:0 |

Page 8

## Clock Replacement Illustration

- Max page table size 4
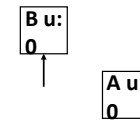
- Invariant: point at oldest page

  – Page B arrives

```
B u:
0
```

## Clock Replacement Illustration

- Max page table size 4

- Invariant: point at oldest page
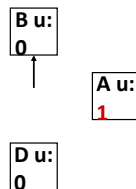
  – Page B arrives
  – Page A arrives
  – Access page A

```
B u:
0
```
```
A u:
0
```

## Clock Replacement Illustration

- Max page table size 4

- Invariant: point at oldest page

  – Page B arrives
  – Page A arrives
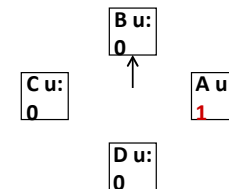  – Access page A
  – Page D arrives

```
B u:
0
```
```
A u:
1
```
```
D u:
0
```

## Clock Replacement Illustration

- Max page table size 4

- Invariant: point at oldest page

  – Page B arrives
  – Page A arrives
  – Access page A
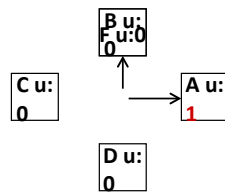  – Page D arrives
  – Page C arrives

```
B u:
0
```
```
C u:
0
```
```
A u:
1
```
```
D u:
0
```

Page 9

## Clock Replacement Illustration

- Max page table size 4

- Invariant: point at oldest page

  – Page B arrives
  – Page A arrives
  – Access page A
  – Page D arrives
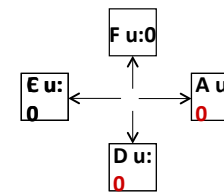  – Page C arrives
  – Page F arrives



**B u:**
**F u:0**
**0**

**C u:**
**0**

**A u:**
**1**

**D u:**
**0**

---

## Clock Replacement Illustration

- Max page table size 4

- Invariant: point at oldest page

  – Page B arrives
  – Page A arrives
  – Access page A
  – Page D arrives
  – Page C arrives
  – Page F arrives
  – Access page D
  – Page E arrives



**F u:0**

**E u:**
**0**

**A u:**
**0**

**D u:**
**0**

---

## N[th] Chance version of Clock Algorithm

- N[th] chance algorithm: Give page N chances
  – OS keeps counter per page: # sweeps
  – On page fault, OS checks use bit:
    » 1⇒clear use and also clear counter (used in last sweep)
    » 0⇒increment counter; if count=N, replace page
  – Means that clock hand has to sweep by N times without page being used before page is replaced
- How do we pick N?
  – Why pick large N? Better approx to LRU
    » If N ~ 1K, really good approximation
  – Why pick small N? More efficient
    » Otherwise might have to look a long way to find free page
- What about dirty pages?
  – Takes extra overhead to replace a dirty page, so give dirty pages an extra chance before replacing?
  – Common approach:
    » Clean pages, use N=1
    » Dirty pages, use N=2 (and write back to disk when N=1)

---

## Clock Algorithms: Details

- Which bits of a PTE entry are useful to us?
  – Use: Set when page is referenced; cleared by clock algorithm
  – Modified: set when page is modified, cleared when page written to disk
  – Valid: ok for program to reference this page
  – Read-only: ok for program to read page, but not modify
    » For example for catching modifications to code pages!
- Do we really need hardware-supported "modified" bit?
  – No.  Can emulate it (BSD Unix) using read-only bit
    » Initially, mark all pages as read-only, even data pages
    » On write, trap to OS. OS sets software "modified" bit, and marks page as read-write.
    » Whenever page comes back in from disk, mark read-only

Page 10

## Clock Algorithms Details (continued)

- Do we really need a hardware-supported "use" bit?

  - No. Can emulate it using "invalid" bit:
    - » Mark all pages as invalid, even if in memory
    - » On read to invalid page, trap to OS
    - » OS sets use bit, and marks page read-only

  - When clock hand passes by, reset use bit and mark page as invalid again

## Summary (1/2)

- TLB is cache on translations
  - Fully associative to reduce conflicts
  - Can be overlapped with cache access

- Demand Paging:
  - Treat memory as cache on disk
  - Cache miss $\Rightarrow$ 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

## Summary (2/2)

- Replacement policies
  - FIFO: Place pages on queue, replace page at end
  - MIN: Replace page that will be used farthest in future
  - LRU: Replace page used farthest in past

- Clock Algorithm: Approximation to LRU
  - Arrange all pages in circular list
  - Sweep through them, marking as not "in use"
  - If page not "in use" for one pass, than can replace

- Second-Chance List algorithm: Yet another approx LRU
  - Divide pages into two groups, one of which is truly LRU and managed on page faults.

Page 11