Welcome back!

# MEMORY: SWAPPING

Shivaram Venkataraman

CS 537, Spring 2023

# ADMINISTRIVIA

Project 3 was due Monday. ?

Project 4: Scheduling new dates: Feb $22^{nd}$ to March $6^{th}$     ~ 2 weeks

$\hookrightarrow$ March 2

Midterm 1: In class midterm, Multiple choice.

No notes / calculators. (We will give a table of powers of 2)

$\rightarrow$ Old exams on Canvas

$\rightarrow$ Discussion: Practice problems

$\quad\quad\quad \hookrightarrow$ Handout

$\rightarrow$ Textbook

Video Playlist

$\hookrightarrow$ Some from past, this year

# OFFICE HOURS

Students

1. One question per student at a time
2. Please be prepared before asking questions
3. The TAs might not be able to fix your problem →
4. Limited time per student → 10 mins

→ Separate branch/ for your problem directory

Run Address Sanitizer

Our

1. Increase number of TAs close to deadline
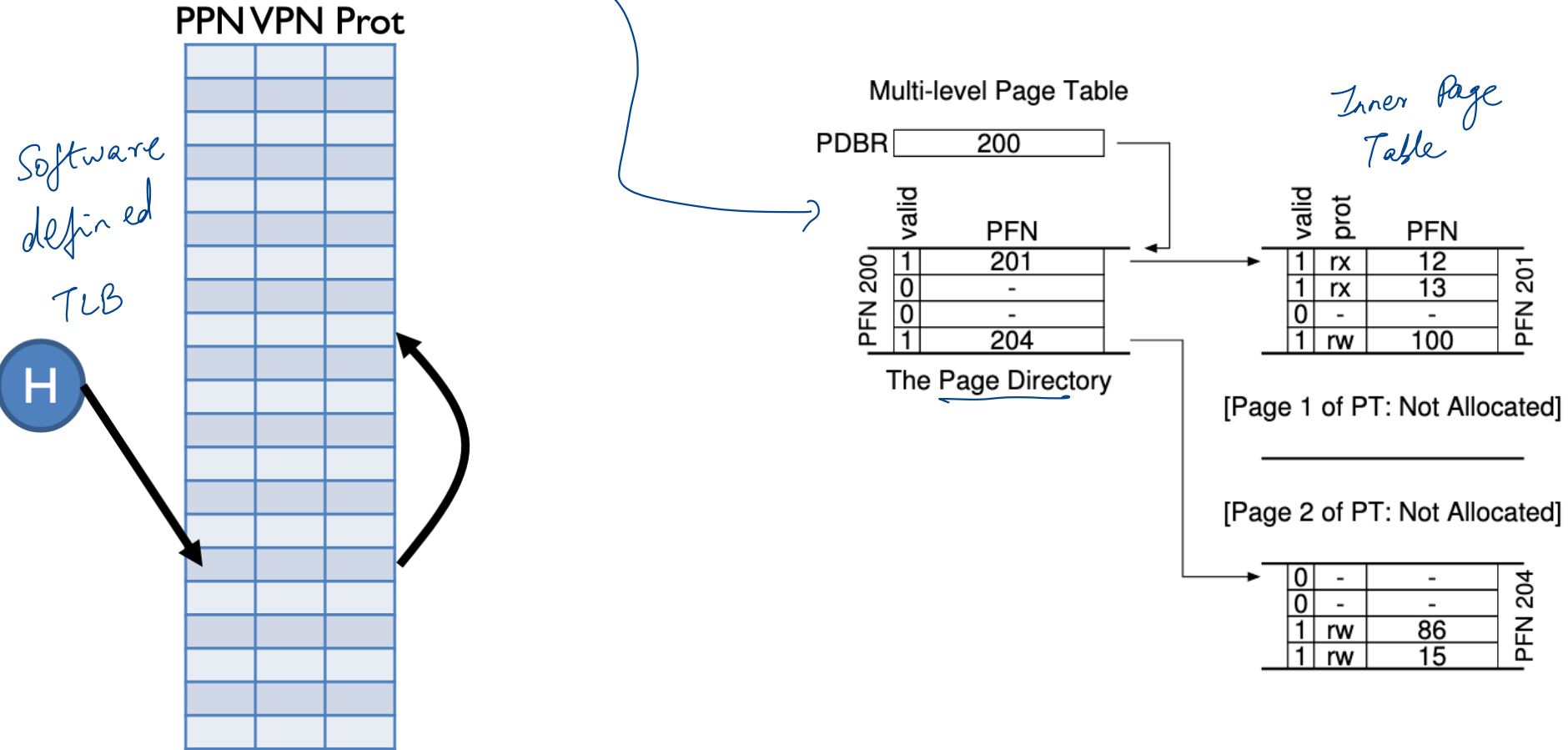2. Study groups
   ↳ Midterm

# AGENDA / LEARNING OUTCOMES

Memory virtualization

How we support virtual mem larger than physical mem?
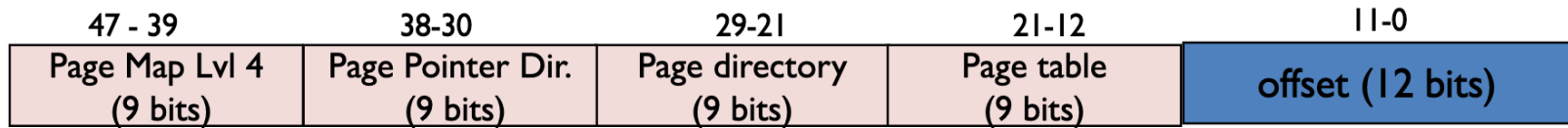
What are mechanisms and policies for this?
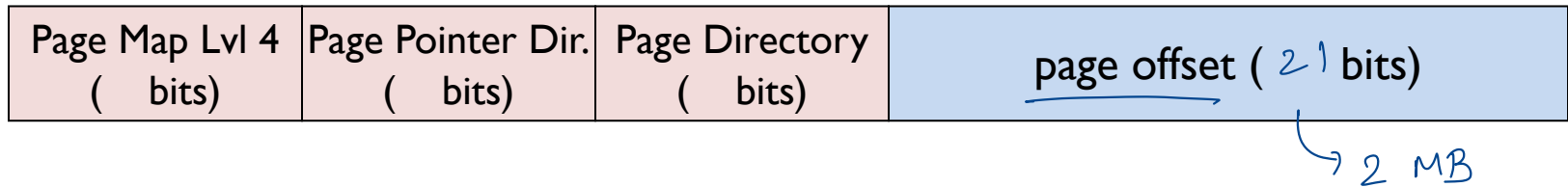
# RECAP

# MULTILEVEL, INVERTED PAGE TABLES

**PPN VPN Prot**

Software defined TLB

H

Multi-level Page Table

Inner Page Table

PDBR | 200 |

| valid | | PFN |
|---|---|---|
| 1 | | 201 |
| 0 | | - |
| 0 | | - |
| 1 | | 204 |

PFN 200

The Page Directory

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

PFN 201

[Page 1 of PT: Not Allocated]

[Page 2 of PT: Not Allocated]

| | | |
|---|---|---|
| 0 | - | - |
| 0 | - | - |
| 1 | rw | 86 |
| 1 | rw | 15 |

PFN 204

# TRANSLATING LARGE PAGES

HugePages saves TLB entries. But how does it affect page translation?

4KB pages: 4 levels → 4 memory accesses

| 47 - 39 | 38-30 | 29-21 | 21-12 | 11-0 |
|---|---|---|---|---|
| Page Map Lvl 4 (9 bits) | Page Pointer Dir. (9 bits) | Page directory (9 bits) | Page table (9 bits) | offset (12 bits) |

2MB pages: 3 levels → 3 memory access on translation

| Page Map Lvl 4 ( bits) | Page Pointer Dir. ( bits) | Page Directory ( bits) | page offset ( 21 bits) |
|---|---|---|---|

↳ 2 MB

# SUMMARY: BETTER PAGE TABLES

Problem:  Simple linear page tables require too much contiguous memory

Many options for efficiently organizing page tables

If OS traps on TLB miss, OS can use any data structure

- Inverted page tables (hashing) $\longrightarrow$ reduce page table size

If Hardware handles TLB miss, page tables must follow specific format

- Multi-level page tables used in x86 architecture
- Each inner page table fits within a page

Large pages can reduce TLB use and number of accesses for translation

Paging
↓
TLB (cache)
↓
Multi Level,
Inverted

# SWAPPING

# MOTIVATION

OS goal: Support processes when not enough physical memory
  - Single process with very large address space
  - Multiple processes with combined address spaces

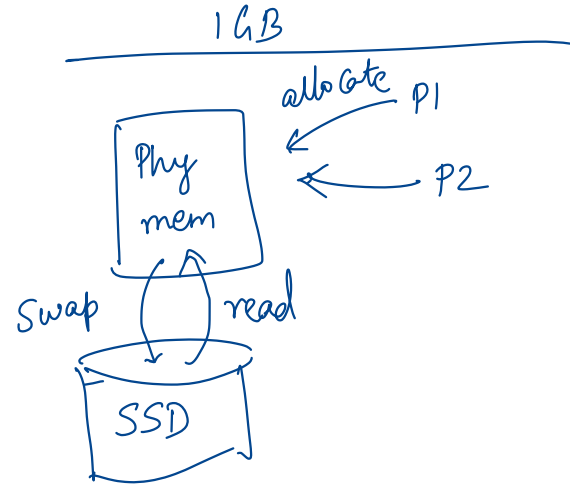User code should be independent of amount of physical memory
  - Correctness, if not performance

Virtual memory: OS provides illusion of more physical memory

Why does this work?
  - Relies on key properties of user processes (workload) and machine architecture (hardware)

Chrome with 100 Tabs each Tab is process using

1 GB

allocate P1

P2

Phy mem

Swap   read

SSD

# WORKLOAD PROPERTIES

*LRU*
*schemes for cache*

Leverage locality of reference within processes

– Spatial: reference memory addresses **near** previously referenced addresses

– Temporal: reference memory addresses that have referenced in the past

*recently*

– Processes spend majority of time in small portion of code
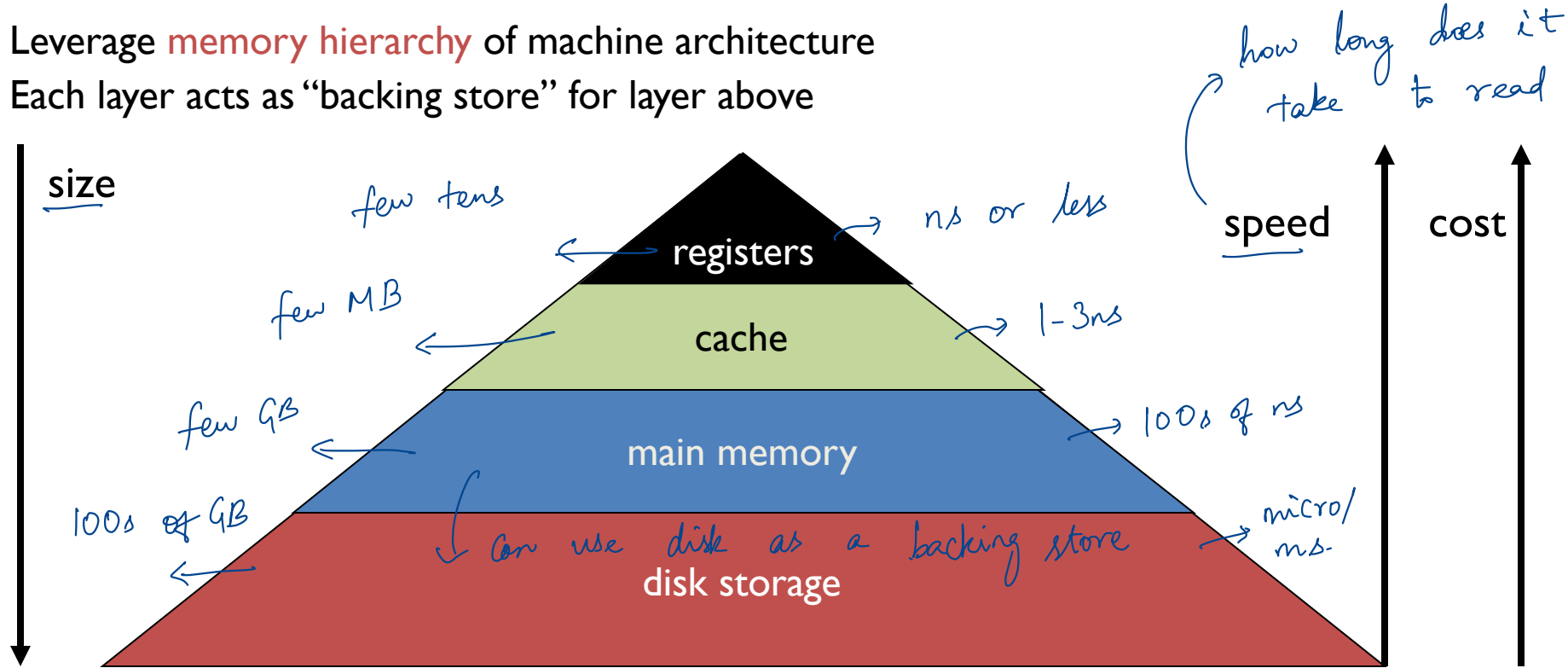
  • Estimate: 90% of time in 10% of code

Implication:

– Process only uses small amount of address space at any moment

– Only small amount of address space must be resident in physical memory

# HARDWARE: MEMORY HIERARCHY

Leverage memory hierarchy of machine architecture
Each layer acts as "backing store" for layer above



size

how long does it take to read

speed        cost

few tens        ns or less

registers

few MB        1-3ns

cache

few GB        100s of ns

main memory

100s of GB        micro/ms.

Can use disk as a backing store

disk storage

# SWAPPING INTUITION

Page Table P1

Idea: OS keeps unreferenced pages on disk
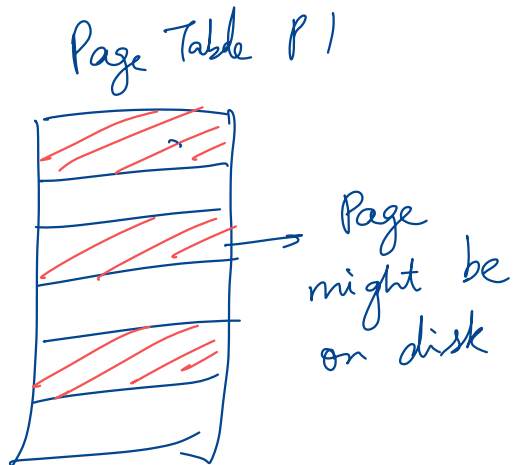- Slower, cheaper backing store than memory

Process can run when not all pages are loaded into main memory

OS and hardware cooperate to make large disk seem like memory
- Same behavior as if all of address space in main memory

Requirements:
- OS must have **mechanism** to identify location of each page in address space → in memory or on disk
- OS must have **policy** to determine which pages live in memory and which on disk
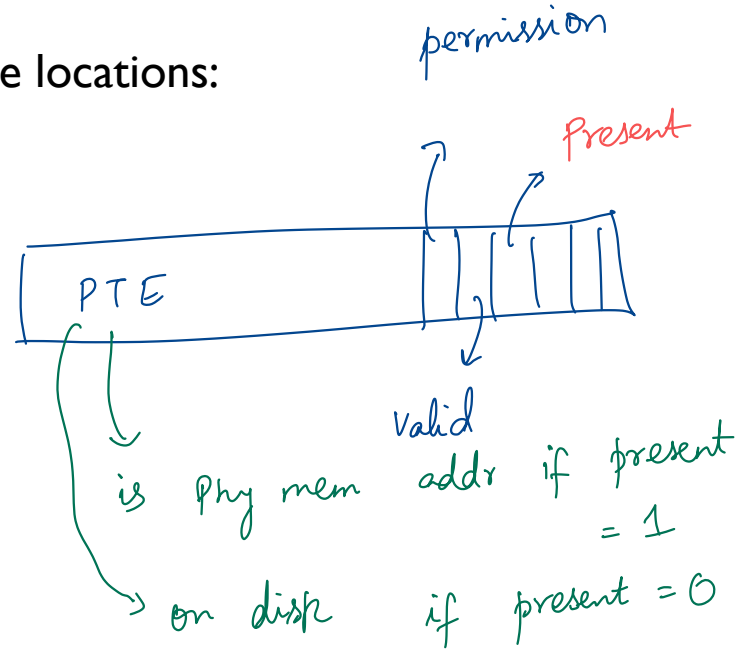
Page might be on disk

# VIRTUAL ADDRESS SPACE MECHANISMS

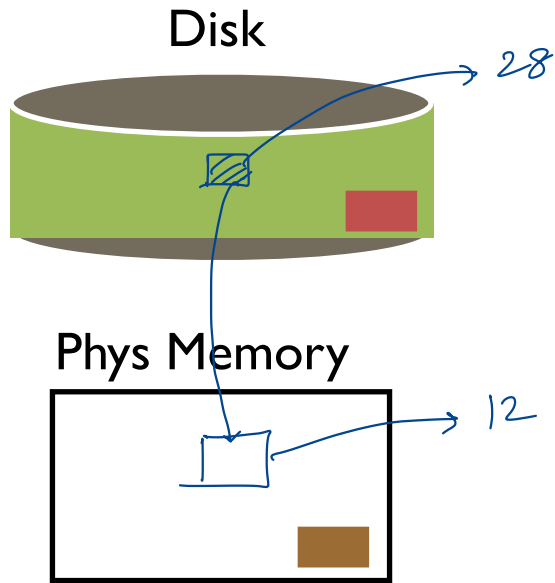Each page in virtual address space maps to one of three locations:

- Physical main memory: Small, fast, expensive
- Disk (backing store): Large, slow, cheap
- Nothing (error): Free

Extend page tables with an extra bit: present

- permissions (r/w), valid, present
- Page in memory: present bit set in PTE
- Page on disk: present bit cleared
  - PTE points to block on disk
  - Causes trap into OS when page is referenced
  - Trap: page fault

permission

Present

PTE

Valid

is Phy mem    addr if present
                        = 1

on disk    if present = 0

Linear Page table

Disk

→ 28

Phys Memory

→ 12

What if access vpn 0xb?

| PFN | valid | prot | present |
|-----|-------|------|---------|
| 10 | 1 | r-x | 1 |
| - | 0 | - | - |
| 23 | 1 | rw- | 0 |
| - | 0 | - | - |
| - | 0 | - | - |
| - | 0 | - | - |
| - | 0 | - | - |
| - | 0 | - | - |
| - | 0 | - | - |
| 12 28 | 1 | rw- | 0 1 |
| 4 | 1 | rw- | 1 |

retry the translation

trap = Page fault

- Read page from block 28
- Store page 12
- Update PTE

# VIRTUAL MEMORY MECHANISMS

First, hardware checks TLB for virtual address
- – if TLB hit, address translation is done; page in physical memory ✓

Else          ...
- – Hardware or OS walk page tables
- – If PTE designates page is present, then page in physical memory ✓
    (i.e., present bit is cleared)

Else

*Policy*
- – Trap into OS (not handled by hardware)
- – OS selects (victim page) in memory to replace
    - • Write victim page out to disk if modified (use dirty bit in PTE)
- – OS reads referenced page from disk into memory
- – Page table is updated, present bit is set
- – Process continues execution

*did the process change the page compared to contents on disk?*

# SWAPPING POLICIES

# SWAPPING POLICIES

Goal: Minimize number of page faults

- Page faults require milliseconds to handle (reading from disk)
- Implication: Plenty of time for OS to make good decision  ⟶ *unlike TLBs*

OS has two decisions
- Page selection
  **When** should a page (or pages) on disk be **brought into** memory?

- Page replacement
  **Which r**esident page (or pages) in memory should be **thrown out** to disk?
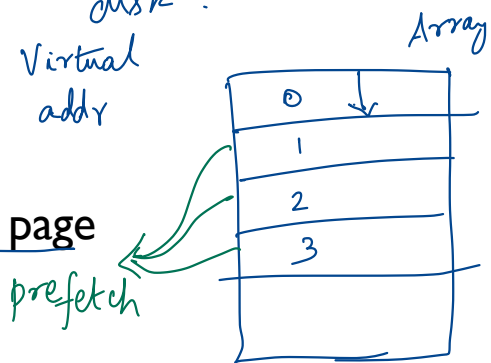  *or evicted*

# PAGE SELECTION

*When should we fetch a page from disk?*

**Demand paging:** Load page only when page fault occurs
- Intuition: Wait until page must absolutely be in memory
- When process starts: No pages are loaded in memory
- Problems: Pay cost of page fault for every newly accessed page

*Virtual addr*

*Array*

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |

*prefetch*

**Prepaging (anticipatory, prefetching):** Load page before referenced
- OS predicts future accesses (oracle) and brings pages into memory early
- Works well for some access patterns (e.g., sequential)

**Hints:** Combine above with user-supplied hints about page references
- User specifies: may need page in future, don't need this page anymore, or sequential access pattern, ...
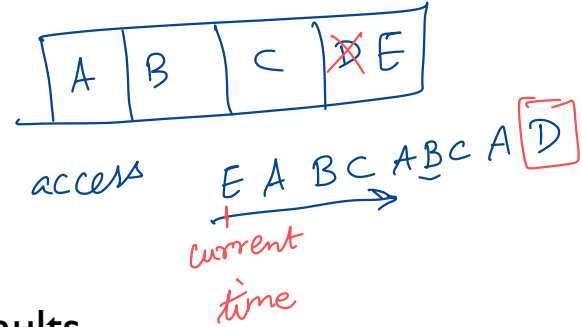- Example: madvise() in Unix

*I will access this page in the future*

# PAGE REPLACEMENT

→ many decades

1960s

Which page in main memory should selected as victim?

- Write out victim page to disk if modified (dirty bit set)
- If victim page is not modified (clean), just discard

| A | B | C | ~~D~~ E |
|---|---|---|---|

stream of page access

E A BC ABC A D
↑
current time

OPT: Replace page not used for longest time in future

- Advantages: Guaranteed to minimize number of page faults
- Disadvantages: Requires that OS predict the future; Not practical, but good for comparison

# PAGE REPLACEMENT

**FIFO**: Replace page that has been in memory the longest      *First In  First Out*
- – Intuition: First referenced long time ago, done with it now
- – Advantages: Fair: All pages receive equal residency; Easy to implement
- – Disadvantage: Some pages may always be needed

   *→ If a page is popular, then it might not stay in memory*

**LRU**: Least-recently-used: Replace page not used for longest time in past
- – Intuition: Use past to predict the future
- – Advantages: With locality, LRU approximates OPT
- – Disadvantages:
  - • Harder to implement, must track which pages have been accessed
  - • Does not handle all workloads well      *sort pages by access time*

# PAGE REPLACEMENT

Three pages of physical memory

Page reference string: DDBBACBDBD

Metric: Miss count

Hit/Miss    OPT        FIFO   6 misses    LRU

| | Hit/Miss | | OPT | | | Hit/Miss | FIFO | | |
|---|---|---|---|---|---|---|---|---|---|---|
| M | D | | D | | | M | D | D | | |
| H | D | | D | | | H | D | D | | |
| M | B | | D | B | | M | B | D | B | |
| H | B | | D | B | | H | B | D | B | |
| M | A | | D | B | X | M | A | (D) | B | A |
| M | C | | D | B | C | M | C | C | B | A |
| H | B | | | | | H | B | C | B | A |
| H | D | | | | | M | D | C | D | A |
| H | B | | | | | M | B | C | D | B |
| H | D | | | | | H | D | | | |

4 misses

# QUIZ 13

https://tinyurl.com/cs537-sp23-quiz13

Page reference string: ABCABDADBCB

Three pages of physical memory

Metric: Miss count

| | | OPT | | | | FIFO = 7 | | | | LRU = 5 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 | ABC | A | B | C | 3 | A | B | C | 3 | A | B | C |
| H | A | | | | | | | | | A | B | C |
| H | B | | | | | | | | | A | B | ~~C~~ |
| 4 | D | A | B | D | ·M | D | B | C | ·M | A | B | D |
| H | A | | | | ·M | D | A | C | H | | | |
| H | D | | | | H | | | | H | | | |
| H | B | | | | ·M | D | A | B | H | A | B | D |
| 5 | C | C | B | D | ·M | C | A | B | ·M | C | B | D |
| H | B | | | | H | | | | H | | | |

# PAGE REPLACEMENT COMPARISON

Add more physical memory, what happens to performance?

LRU, OPT:

- Guaranteed to have fewer (or same number of) page faults
- Smaller memory sizes are guaranteed to contain a subset of larger memory sizes
- Stack property: smaller cache always subset of bigger

FIFO:

- Usually have fewer page faults
- Belady's anomaly: May actually have more page faults!

*access pattern → 3 pages*

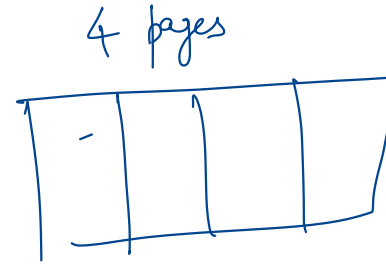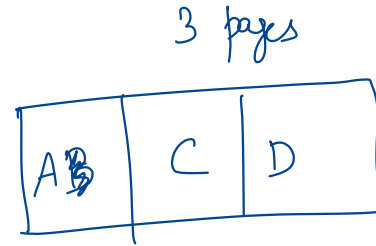*5 misses*

*6 pages of mem.*

*for some access patterns*

# FIFO PERFORMANCE MAY DECREASE!

Consider access stream: ABCDABEABCDE

Physical memory size: 3 pages vs. 4 pages

How many misses with FIFO?

Anamoly

3 pages

| A B | C | D |
|-----|---|---|

4 pages

| - | | | |
|---|---|---|---|

3 M    for    A, B, C

M    for    D

M    for    A

M    for    B

M    for    E

# IMPLEMENTING LRU

Software Perfect LRU
- OS maintains ordered list of physical pages by reference time
- When page is referenced: Move page to front of list
- When need victim: Pick page at back of list
- Trade-off: Slow on memory reference, fast on replacement

Hardware Perfect LRU
- Associate timestamp register with each page
- When page is referenced: Store system clock in register
- When need victim: Scan through registers to find oldest clock
- Trade-off: Fast on memory reference, slow on replacement (especially as size of memory grows)

In practice
LRU is an approximation anyway, so approximate more?

*sorted linked list based on reference*
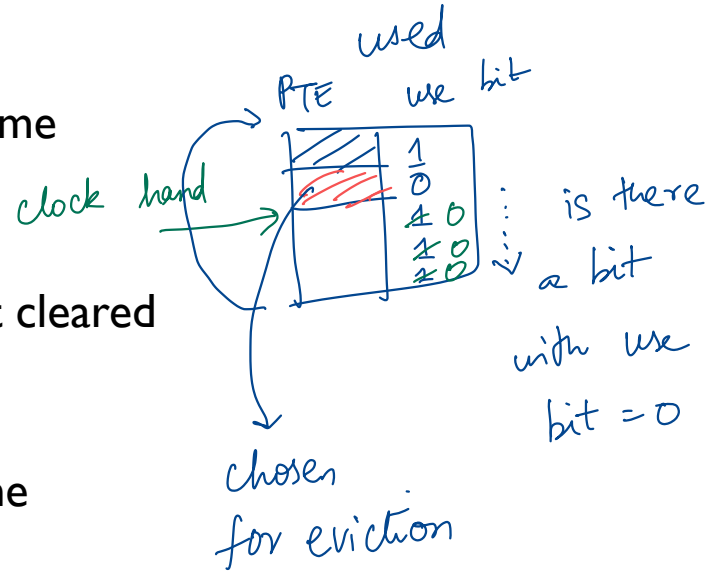
*expensive large num of pages*

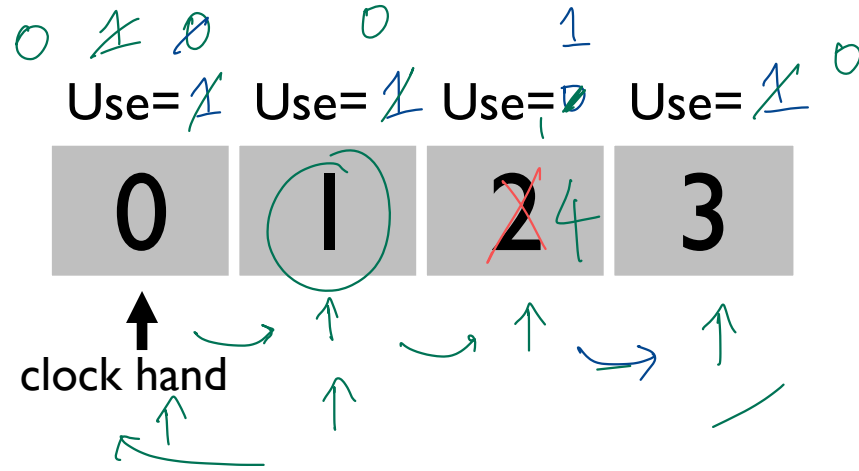# CLOCK ALGORITHM

some old pages not least recently used

Hardware
- Keep use (or reference) bit for each page frame
- When page is referenced: set use bit

Operating System
- Page replacement: Look for page with use bit cleared (has not been referenced for awhile)
- Implementation:
  - Keep pointer to last examined page frame
  - Traverse pages in circular buffer
  - Clear use bits as search
  - Stop when find page with already cleared use bit, replace this page

PTE   use bit

clock hand

1
0
1 0
1 0
1 0

is there a bit with use bit = 0

chosen for eviction

# CLOCK: LOOK FOR A PAGE

0  1  0        0        1

Use=1  Use=1  Use=0  Use=1  0

Physical Mem:



| 0 | 1 | 2 4 | 3 |

↑
clock hand

Use = 1,1,0,1 to begin

Evict a page    bring in 4
    ↳ select   2 for eviction

→ Page 0 is accesed

Page 5  Bring in
    ↳ select  1 for
            eviction

# CLOCK EXTENSIONS

Replace multiple pages at once
- – Intuition: Expensive to run replacement algorithm and to write single block to disk
- – Find multiple victims each time and track free list


Use dirty bit to give preference to dirty pages
- – Intuition: More expensive to replace dirty pages
  Dirty pages must be written to disk, clean pages do not
- – Replace pages that have use bit and dirty bit cleared

# SUMMARY: VIRTUAL MEMORY

Abstraction: Virtual address space with code, heap, stack

Address translation

- Contiguous memory: base, bounds, segmentation

- Using fixed sizes pages with page tables

Challenges with paging

- Extra memory references: avoid with TLB

- Page table size: avoid with multi-level paging, inverted page tables etc.

Larger address spaces: Swapping mechanisms, policies (LRU, Clock)

# NEXT STEPS

Next class: New module on Concurrency!