# 18-344 Recitation 4

09/26/2025

#### Outline

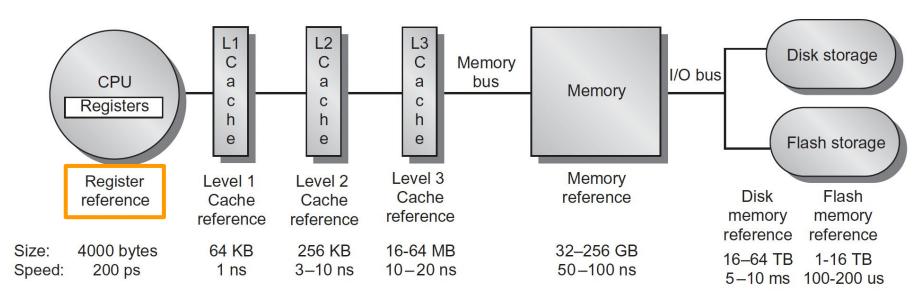
- 1. Logistics
- 2. Review
- 3. Homework 1 postmortem

#### Logistics

- Lab 0 and Homework 1 grades have been released
- Lab 1 released
  - Code freeze: September 28 (Sunday) (in 2 days)
    - Submit (push a commit) on GitHub Classroom
  - Report submission: October 1 (Wednesday) (in 5 days)
    - Submit on Gradescope
- Homework 3 released, due September 29 (Monday) (in 3 days)
  - Covers materials from lectures 7-9
    - Mostly caching
- Homework 2 grades to be released around this weekend/early next week

# Review

## The memory hierarchy



Memory hierarchy for server

Computer Architecture: A Quantitative Approach 6th Edition, Hennesy and Patterson

## Memory access time

Generally,

Average memory access time = Hit time + Miss rate  $\times$  Miss penalty

#### Caches' 3 C's

- Compulsory miss: The first access to a block address is always a miss
- Capacity miss: Cache is too small to hold everything long enough so they can be reused
  - Think of it as a non-compulsory miss that would occur in an ideal fully-associative cache
- Conflict miss: Miss to a previously visited block that was evicted due to another block address being mapped to the same cache block
  - If a miss isn't a compulsory or a capacity miss, it's a conflict miss.

- Round-robin
- Least-recently used (LRU)

- Round-robin
- Least-recently used (LRU)

```
accessCacheLRU(access a) {
  for each block in cache, b:
    if b != a.block:
      LRU Age[b]++
  LRU Age[b] = 0
findBlockLRU() {
  return argmax (b, LRU Age)
```

- Round-robin
- Least-recently used (LRU)
  - Expensive to implement
    - Area & Power cost to store ages
    - Time & Energy cost to update ages and identify the block to evict

- Round-robin
- Least-recently used (LRU)
  - Expensive to implement
    - Area & Power cost to store ages
    - Time & Energy cost to update ages and identify the block to evict
- Bit-pseudo-least-recently used (BPLRU)
  - Evict a block that was definitely not most recently used, decent approximation of LRU

- Round-robin
- Least-recently used (LRU)
  - Expensive to implement
    - Area & Power cost to store ages
    - Time & Energy cost to update ages ar
- Bit-pseudo-least-recently used (BPLR)
  - Evict a block that was definitely not most re

```
accessCachePLRU(access a) {
  MRU Bit[a.block] = 1
  if ++MRU BitSum == setSize:
    for each block in cache, b:
      MRU Bit[b] = 0
    MRU BitSum = 0
findBlockLRU() {
  for i in 0..setSize:
    if !MRU Bit[i]:
      return block(i);
```

- Round-robin
- Least-recently used (LRU)
  - Expensive to implement
    - Area & Power cost to store ages
    - Time & Energy cost to update ages and identify the block to evict
- Bit-pseudo-least-recently used (BPLRU)
  - Evict a block that was definitely not most recently used, decent approximation of LRU
  - Not as expensive as LRU

- Round-robin
- Least-recently used (LRU)
  - Expensive to implement
    - Area & Power cost to store ages
    - Time & Energy cost to update ages and identify the block to evict
- Bit-pseudo-least-recently used (BPLRU)
  - Evict a block that was definitely not most recently used, decent approximation of LRU
  - Not as expensive as LRU
- Belady's MIN algorithm (MIN)
  - Optimally evict a block that has the longest reuse distance

- Round-robin
- Least-recently used (LRU)
  - Expensive to implement
    - Area & Power cost
    - Time & Energy cos
- Bit-pseudo-least-recent
  - Evict a block that was de
  - Not as expensive as LRI
- Belady's MIN algorithm
  - Optimally evict a block th

```
findBlockMIN() {
    0://init reuse distances
    1: for each block in cache, b:
          RD[b] = 0; RD done[b] = false;
    3://look forward in the execution trace
    4: for each access, a, forward in execution trace:
    5://increment reuse distance for each block not already seen
    6:
          for each block in cache, b:
    7:
               if RD done[b] == false:
    8:
                  RD[b]++;
    9:
          RD done[a.block] = true
   10://MIN finds the block with maximum RD
   11:return argmax(b,RD[b])
```

- Round-robin
- Least-recently used (LR<u>U)</u>
  - Expensive to implement
    - Area & Power cost
    - Time & Energy cos
- Bit-pseudo-least-recentl
  - Evict a block that was de
  - Not as expensive as LRU
- Belady's MIN algorithm
  - Optimally evict a block th

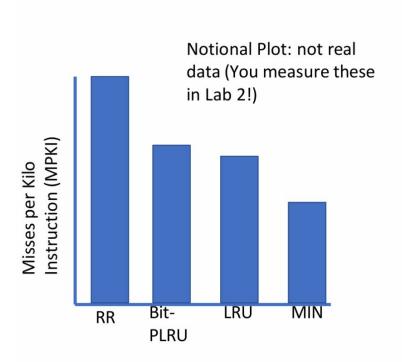
```
findBlockMIN() {
    0://init reuse distances
    1:for each block in cache, b:
    2:    RD[b] = 0; RD_done[b] = false;
    3://look forward in the execution trace
    4:for each access, a, forward in execution trace:
    5://increment reuse distance for each block not already seen
    6:    for each block in cache, b:
    7:        if RD_done[b] == false:
    8:        RD[b]++;
    9:    RD_done[a.block] = true
    10://MIN finds the block with maximum RD
    11:return argmax(b,RD[b])
```

- Round-robin
- Least-recently used (LRU)
  - Expensive to implement
    - Area & Power cost to store ages
    - Time & Energy cost to update ages and identify the block to evict
- Bit-pseudo-least-recently used (BPLRU)
  - Evict a block that was definitely not most recently used, decent approximation of LRU
  - Not as expensive as LRU
- Belady's MIN algorithm (MIN)
  - Optimally evict a block that has the longest reuse distance
  - Not realistic to implement

- Round-robin
- Least-recently used (LRU)
  - Expensive to implement
    - Area & Power cost to store ages
    - Time & Energy cost to update ages and identify the block to evict
- Bit-pseudo-least-recently used (BPLRU)
  - Evict a block that was definitely not most recently used, decent approximation of LRU
  - Not as expensive as LRU
- Belady's MIN algorithm (MIN)
  - Optimally evict a block that has the longest reuse distance
  - Not realistic to implement

Capacity miss. Think of it as a non-compulsory miss that would occur in an *ideal* fully-associative cache.

#### Comparing cache replacement policies



**RR:** log(set size) bits per set to track next to evict, no action on access

**Bit-PLRU:** 1 MRU bit per block + log(set size) bits per set (or equivalent logic) to detect all set,
Clear bits on access if all bits set

**LRU:** 1 age per block + logic to track max. Update (set size - 1) ages on any access

**MIN:** unimplementable, requires future knowledge of execution trace.

## Caching optimizations

#### Victim cache

(Usually) small, fully-associative cache that tracks evicted blocks.

#### Stream buffer

 A buffer that prefetches successive blocks from a lower level cache to a higher level cache upon a miss.

#### Write buffer

- Single-cycle operation to buffer a write request before it is issued to the cache.
- Ordering challenges.

#### Scratchpad

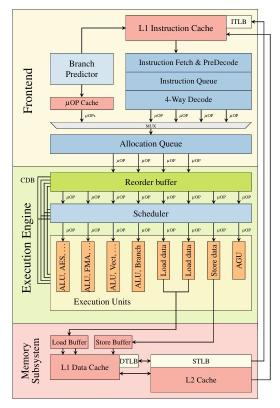
Software controlled memory with explicit, scratch-pad-private physical memory space.

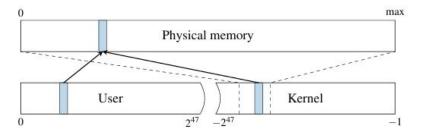


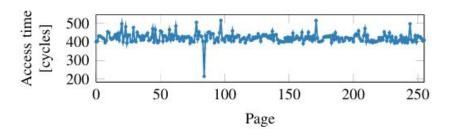


and — we'll talk about these briefly

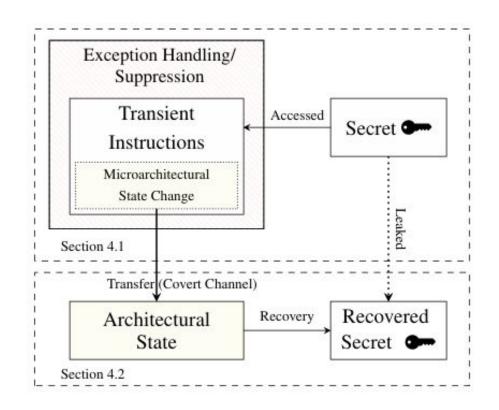




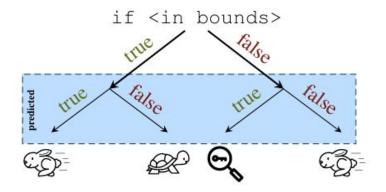


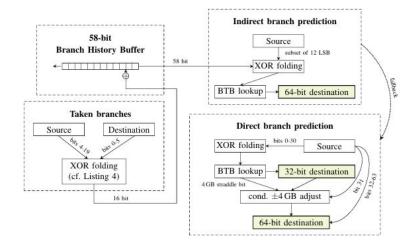


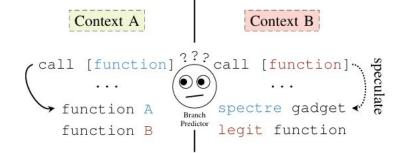












Homework 1 postmortem

- all of you did a pretty good job

#### Common questions

Q1.2

3 Points

Your friend proposes a operations by 30%. Cal

Answer in nanosecond

#### Speedup w/ Example

"You friend proposes an optimization which would speed up *all* operations by 30%"

$$X = 1.3 = \frac{T_{\text{original}}}{T_{\text{optimized}}}$$