

308-273

Caches, Part II

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Review

- We would like to have the capacity of disk at the speed of the processor: unfortunately this is not feasible.
- So we create a memory hierarchy:
 - each successive level contains “most used” lower level
 - exploits **temporal** and **spatial** locality
 - do the common case fast, worry less about the exceptions (design principle of MIPS)
- Locality of reference is a Big Idea

Big Idea Review (1/2)

- Mechanism for transparent movement of data among levels of a storage hierarchy
 - set of address/value bindings
 - address \Rightarrow index to set of candidates
 - compare desired address with tag
 - service hit or miss
 - load n on miss
- <https://eduassistpro.github.io/>

address:	tag	Add WeChat edu_assist_pro	offset
00000000	000000000000	000000000001	1100

The diagram shows a 4x4 grid representing memory. The columns are labeled at the top: **Tag**, **0x0-3**, **0x4-7**, **0x8-b**, and **0xc-f**. The first column (Tag) has a red '1' in the second row and a red '0' in the first row. A red arrow points from the word 'Valid' to the '0' in the first row of the Tag column. Another red arrow points from the top right to the '0xc-f' header. The cell containing 'd' in the first row of the '0xc-f' column is circled in red.

	Tag	0x0-3	0x4-7	0x8-b	0xc-f
	0				d
1		a	b	c	

Outline

- **Block Size Tradeoff**
- **Types of Cache Misses**
- **Fully Associative Cache**
Assignment Project Exam Help
- **N-Way Ass** <https://eduassistpro.github.io/>
- **Block Replacement** Add WeChat edu_assist_pro
- **Multilevel Caches (if time)**
- **Cache write policy (if time)**

Block Size Tradeoff (1/3)

◦ Benefits of Larger Block Size

- **Spatial Locality**: if we access a given word, we're likely to access other nearby words soon (Another Big Idea)
- **Very applicable** - Program
Concept: <https://eduassistpro.github.io/>
instruction, it's likely 'll execute the next few as well
- Works nicely in sequential array accesses too

Block Size Tradeoff (2/3)

◦ Drawbacks of Larger Block Size

- Larger block size means **larger miss penalty**
 - on a miss, takes longer time to load a new block from next level
- If block size is close to cache size, then the miss rate is high
- Result: miss rate goes up

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

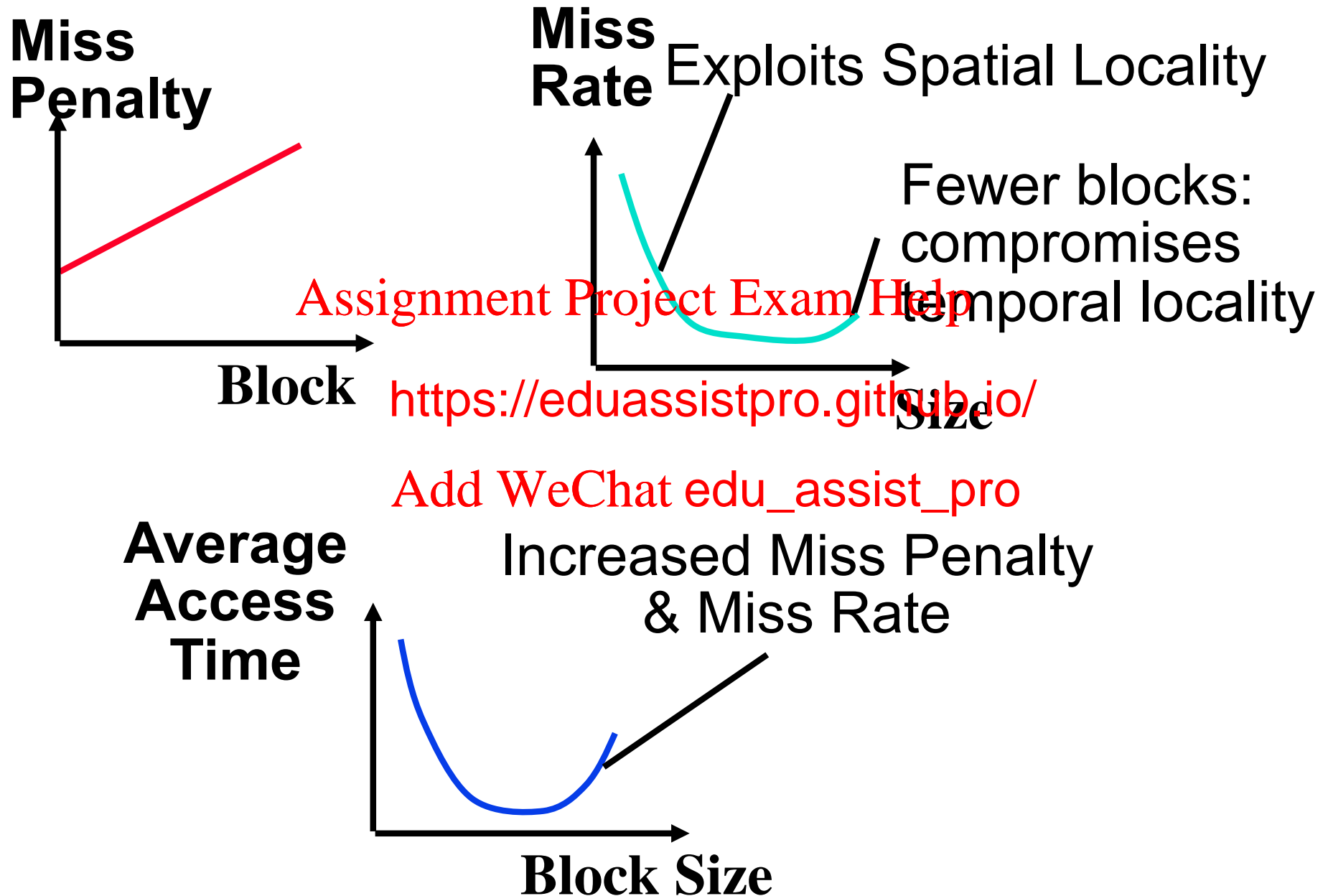
◦ In general, minimize **Average Access Time**

$$= \text{Hit Time} + \text{Miss Penalty} \times \text{Miss Rate}$$

Block Size Tradeoff (3/3)

- **Hit Time** = time to find and retrieve data from current level cache
- **Miss Penalty** = average time to retrieve data on a current level miss (includes the possible number of successive levels of memory)
 Assignment Project Exam Help
 <https://eduassistpro.github.io/>
 Add WeChat edu_assist_pro
- **Hit Rate** = % of requests that are found in current level cache
- **Miss Rate** = $1 - \text{Hit Rate}$

Block Size Tradeoff Conclusions



Types of Cache Misses (1/2)

◦ Compulsory Misses

- occur when a program is first started
 - cache does not contain any of that program's data yet, so misses are bound to occur
 - can't be avoided easily
- <https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
- Assignment Project Exam Help
- on't focus on these in this course

Types of Cache Misses (2/2)

◦ Conflict Misses

- miss that occurs because two distinct memory addresses map to the same cache location
- two blocks to map to the same location, each overwriting the other
- big problem in direct-mapped caches
- how do we lessen the effect of these?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Dealing with Conflict Misses

- **Solution 1: Make the cache size bigger**
 - relatively expensive
- **Solution 2: Multiple distinct blocks can fit in the same Cache Index?**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Fully Associative Cache (1/3)

◦ Memory address fields:

- Tag: same as before
- Offset: same as before
- Index: non-existent

Assignment Project Exam Help

◦ What does

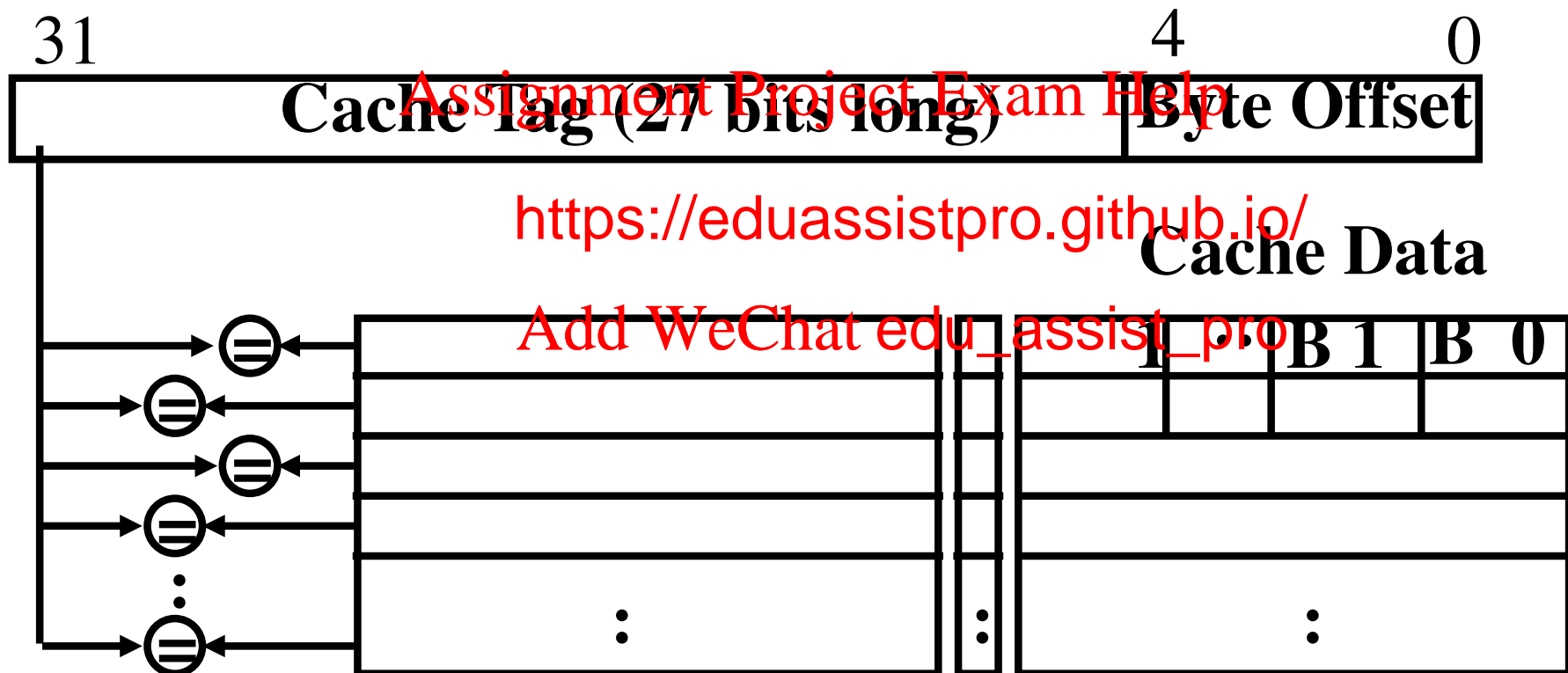
<https://eduassistpro.github.io/>

- any block can go anywhere in the cache
- must compare with all tags in entire cache to see if data is there

Add WeChat edu_assist_pro

Fully Associative Cache (2/3)

- ° Fully Associative Cache (e.g., 32 B block)
 - compare tags in parallel



Fully Associative Cache (3/3)

◦ Benefit of Fully Assoc Cache

- no Conflict Misses (since data can go anywhere)

◦ Drawback

- need hardware comparators: very expensive

◦ Small fully associative cache may be feasible

Third Type of Cache Miss

◦ Capacity Misses

- miss that occurs because the cache has a limited size
- miss that would not occur if we increase the size

<https://eduassistpro.github.io/>

◦ This is the primary miss for Fully Associate caches

Add WeChat edu_assist_pro

N-Way Set Associative Cache (1/4)

◦ Memory address fields:

- Tag: same as before
- Offset: same as before
- Index: position of the correct “row”
(called a <https://eduassistpro.github.io/>

◦ So what's the difference?

- each set contains multiple blocks
- once we've found correct set, must compare with all tags in that set to find our data

N-Way Set Associative Cache (2/4)

◦ Summary:

- cache is direct-mapped with respect to sets
- each set is fully associative
- basically <https://eduassistpro.github.io/> mapped caches, each of which is full associative. Each has its own valid bit

N-Way Set Associative Cache (3/4)

◦ Given memory address:

- Find correct set using Index value.
- Compare Tag with all Tag values in the determined set.
- If a match, otherwise a miss.
- Finally, use the offsets as usual to find the desired data within the desired block.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

N-Way Set Associative Cache (4/4)

◦ What's so great about this?

- even a 2-way set assoc cache avoids a lot of conflict misses
- hardware cost isn't that bad: only need N comparators

Assignment Project Exam Help

<https://eduassistpro.github.io/>

◦ In fact, for blocks,

Add WeChat edu_assist_pro

- it's Direct-Mapped if 1 block per set
- it's Fully Assoc if it's M-way set assoc (M blocks per set)
- so these two are just special cases of the more general set associative design

Block Replacement Policy (1/2)

- **Direct-Mapped Cache:** index completely specifies which position a block can go in on a miss
- **N-Way Set Assoc ($N > 1$):** index specifies a position, an occupy any position on a miss
<https://eduassistpro.github.io/>
- **Fully Associative:** block can be written into any position (there is no index)
Add WeChat edu_assist_pro
- **Question:** if we have the choice, where should we write an incoming block?

Block Replacement Policy (2/2)

◦ **Solution!**

◦ **If there are any locations with valid bit off (empty), then usually write the new block into**

Assignment Project Exam Help

◦ **If all possible locations already have a valid block, we must**
replacement policy by which we determine which block gets “cached out” on a miss.

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Block Replacement Policy: LRU

◦ LRU (Least Recently Used)

- Idea: cache out block which has been accessed (read or write) least recently
- Pro: temporal locality Assignment Practice Exam Help \Rightarrow recent past use implies likelihood, this is a very effective fact, this is a <https://eduassistpro.github.io/>
- Con: with 2-way set Add WeChat edu_assist_pro easy to keep track (one LRU bit); with 4-way or greater, requires complicated hardware and much time to keep track of this

Block Replacement Example

- We have a 2-way set associative cache with a four word *total* capacity and one word blocks. We perform the following word accesses (ignore bytes for this problem):

Assignment Project Exam Help

0, 2, 0, <https://eduassistpro.github.io/>

How many hits and misses will there be for the LRU block replacement policy?

Hint: treat addresses as TAG + INDEX

Block Replacement Example: LRU

◦ Addresses 0, 2, 0, 1, 4, 0, ...

• 0: miss, bring into set 0 (loc 0)

• 2: miss, bring into set 0 (loc 1)

Assignment Project Exam Help
hit

<https://eduassistpro.github.io/>

• 1: miss, bring into set 1 (loc 0)

• 4: miss, bring into set 0 (loc 1, replace 2)

• 0: hit

	loc 0	loc 1
set 0	0	<i>lru</i>
set 1		
set 0	<i>lru</i> 0	2
set 1		
set 0	0	<i>lru</i> 2
set 1		
set 0	0	<i>lru</i> 2
set 1	1	<i>lru</i>
set 0	<i>lru</i> 0	4
set 1	1	<i>lru</i>
set 0	0	<i>lru</i> 4
set 1	1	<i>lru</i>

Ways to reduce miss rate

◦ Larger cache

- limited by cost and technology
- hit time of first level cache $<$ cycle time

Assignment Project Exam Help

◦ More places to put each block of memory

- fully-associative
 - any block any line
- k-way set associated
 - k places for each block
 - direct map: $k=1$

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Big Idea

- How do we choose between options of associativity, block size, replacement policy?
- Design against a performance model
 - **Minimize** <https://eduassistpro.github.io/>
= Hit Time + Miss x Miss Rate
 - **influenced by technology and program behavior**

Example

◦ Assume

- Hit Time = 1 cycle

- Miss rate = 5%

- Miss pen

◦ Avg mem a

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

05 x 20

ycle

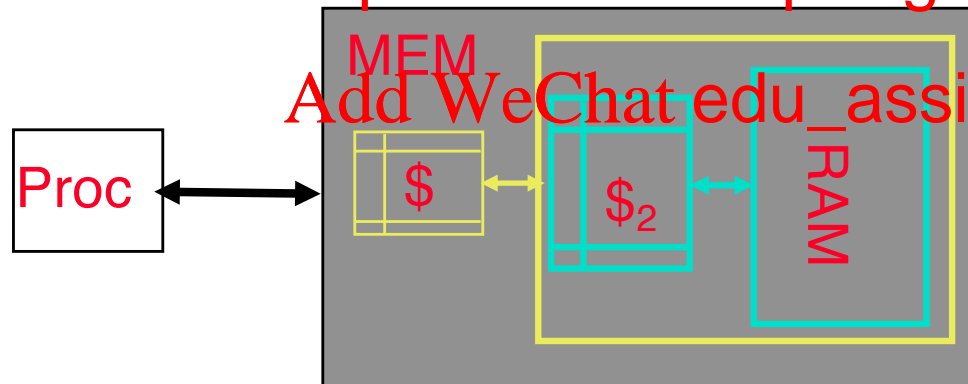
Improving Miss Penalty

- When caches first became popular, Miss Penalty ~ 10 processor clock cycles
- Today 1000 MHz Processor (1 ns per clock cycle) and 100 ns to go to DRAM
 $\Rightarrow 100$ processor clock cycles!

Assignment Project Exam Help

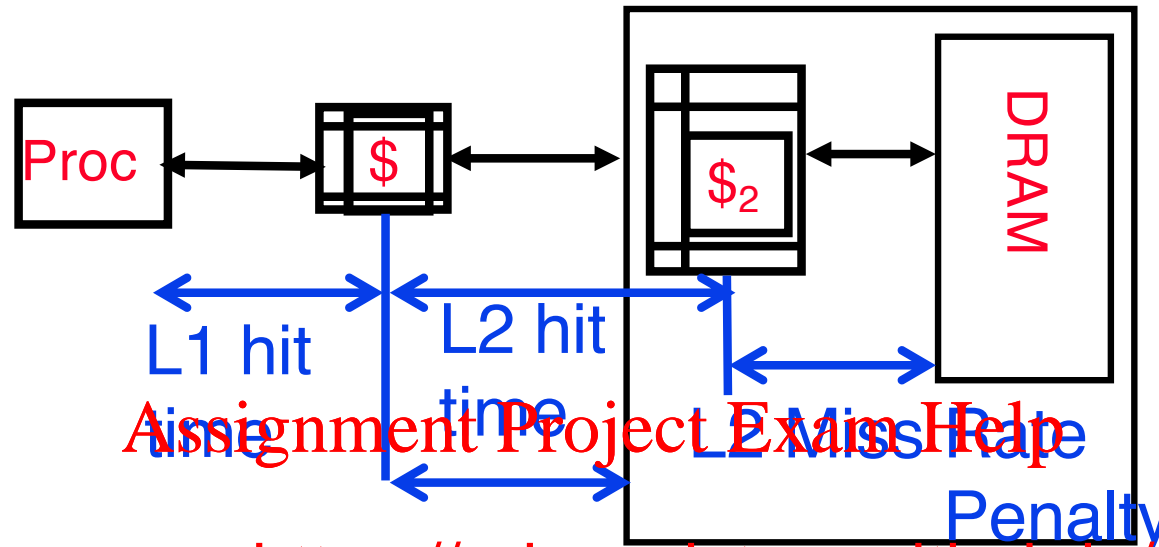
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Solution: another cache between memory and the processor cache: Second Level (L2) Cache

Analyzing Multi-level cache hierarchy



<https://eduassistpro.github.io/>

L1 Miss P

Add WeChat edu_assist_pro

Avg Mem Access Time =

L1 Hit Time + L1 Miss Rate * L1 Miss Penalty

L1 Miss Penalty = L2 Hit Time + L2 Miss Rate

*** L2 Miss Penalty**

Avg Mem Access Time =

**L1 Hit Time + L1 Miss Rate * (L2 Hit Time +
L2 Miss Rate * L2 Miss Penalty)**

Typical Scale

◦ L1

- size: tens of KB
- hit time: complete in one clock cycle
- miss rates: 1-5%

Assignment Project Exam Help

◦ L2:

- size: hundreds of KB
- hit time: few clock cycles
- miss rates: 10-20%

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

◦ L2 miss rate is fraction of L1 misses that also miss in L2

- why so high?

Example: without L2 cache

◦ Assume

- L1 Hit Time = 1 cycle

- L1 Miss rate = 5%

- L1 Miss Penalty = 6 cycles

◦ Avg mem access time

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

5 x 100

6 cycles

Example with L2 cache

◦ Assume

- L1 Hit Time = 1 cycle

- L1 Miss rate = 5%

- L2 Hit Time

- L2 Miss rate

- L2 Miss Penalty = 10

◦ L1 miss penalty = $5 + 0.15 * 100 = 20$

◦ Avg mem access time = $1 + 0.05 * 20$ = 2 cycle

◦ 3x faster with L2 cache

What to do on a write hit?

◦ Write-through

- update the word in cache block and corresponding word in memory

◦ Write-back

- update word
- allow memory word

=> add 'dirty' bit to each line indicating that memory needs to be updated when block is replaced

=> OS flushes cache before I/O !!!

◦ Performance trade-offs?

“And in conclusion...” (1/2)

- Caches are NOT mandatory:
 - Processor performs arithmetic
 - Memory stores data
 - Caches simply make data transfers go faster
- Each level <https://eduassistpro.github.io/> archy is just a subset of next high Add WeChat edu_assist_pro
- Caches speed up due to **temporal locality**: store data used recently
- Block size > 1 word speeds up due to **spatial locality**: store words adjacent to the ones used recently

“And in conclusion...” (2/2)

◦ **Cache design choices:**

- **size of cache: speed v. capacity**
- **direct-mapped v. associative**
- **for N-way set assoc. choice of N**
- **block repl**
- **2nd level cache?**
- **Write through v. write back?**

◦ **Use performance model to pick between choices, depending on programs, technology, budget, ...**