



Microprocessor & Computer Architecture (μ pCA)

UE19CS252

Dr. D. C. Kiran

Department of
Computer Science and Engineering

Microprocessor & Computer Architecture (μ pCA)

Unit 4: Cache Optimization

Dr. D. C. Kiran

Department of Computer Science and Engineering

Microprocessor & Computer Architecture (μpCA)

Syllabus

~~Unit 1: Basic Processor Architecture and Design~~

~~Unit 2: Pipelined Processor and Design~~

~~Unit 3: Memory~~

Unit 4: Input/Output Device Design

~~3rd~~

~~Introduction to Cache Optimization~~

~~Reduce Miss Rate~~

Reduce Miss Penalty

4th Optimization: Multilevel Caches to Reduce Miss Penalty

5th Optimization: Giving priority to Read misses over Write misses to
reduce miss penalty

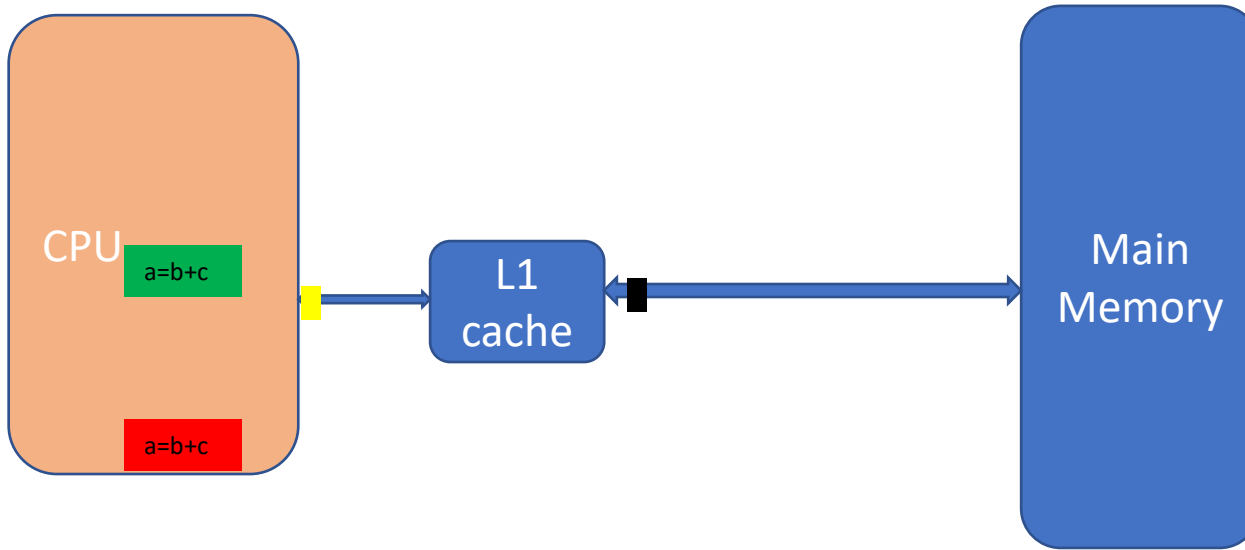
Unit 5: Advanced Architecture



Microprocessor & Computer Architecture (μpCA)

What is the Problem?

$$\text{AMAT} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

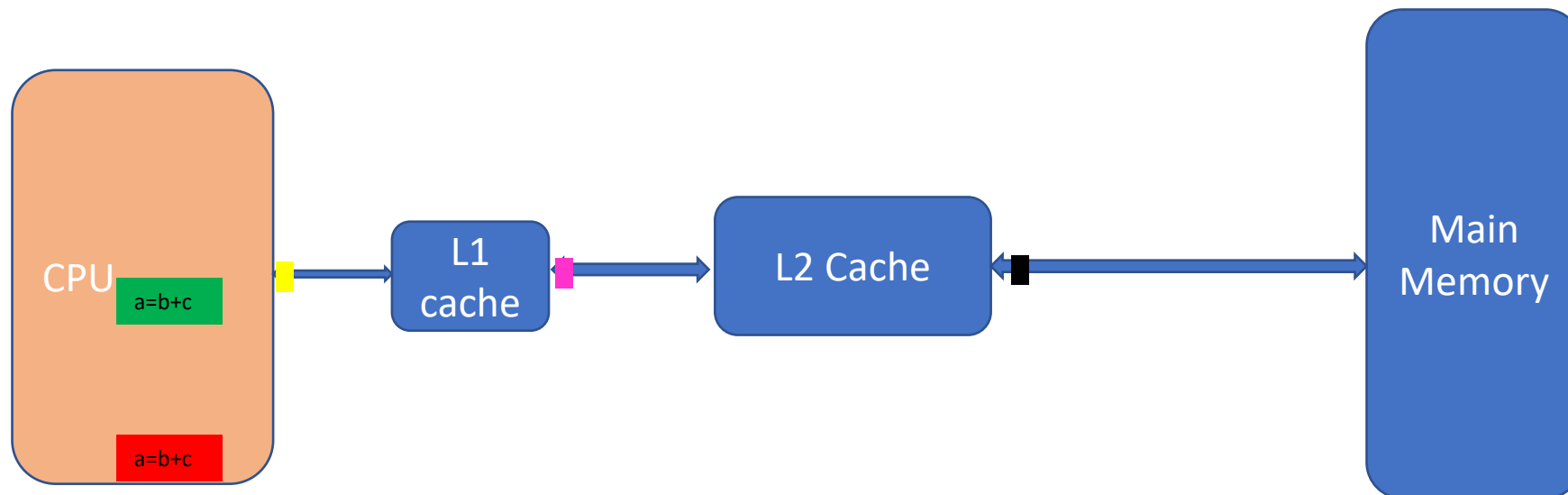


- **Reducing Miss penalty is also important similar to Reducing miss rate.**
- The performance gap between processor & memory raises a question:
 - Should I make the cache faster to keep the pace with the speed of the processor? Or
 - Make the cache larger to overcome the widening gap between the processor & the main memory?

Microprocessor & Computer Architecture (μpCA)

What is the Problem?

- Answer for these questions is to do both.
 - Adding another level of cache between memory & original cache simplifies the decision.
- First level cache can be small enough to match the clock cycle time of the processor.
- Second level cache can be large enough to capture many accesses that would go to main memory, thus reducing miss penalty.



AMAT? For two Level Cache

- The original formula is:

$$\text{AMAT} = \text{Hit Time}_{L1} + \text{Miss Rate}_{L1} \times \text{Miss Penalty}_{L1}$$

Miss in Level 1 Cache, leads to access data from Level 2 cache

$$\text{Miss Penalty}_{L1} = \text{Hit Time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2}$$

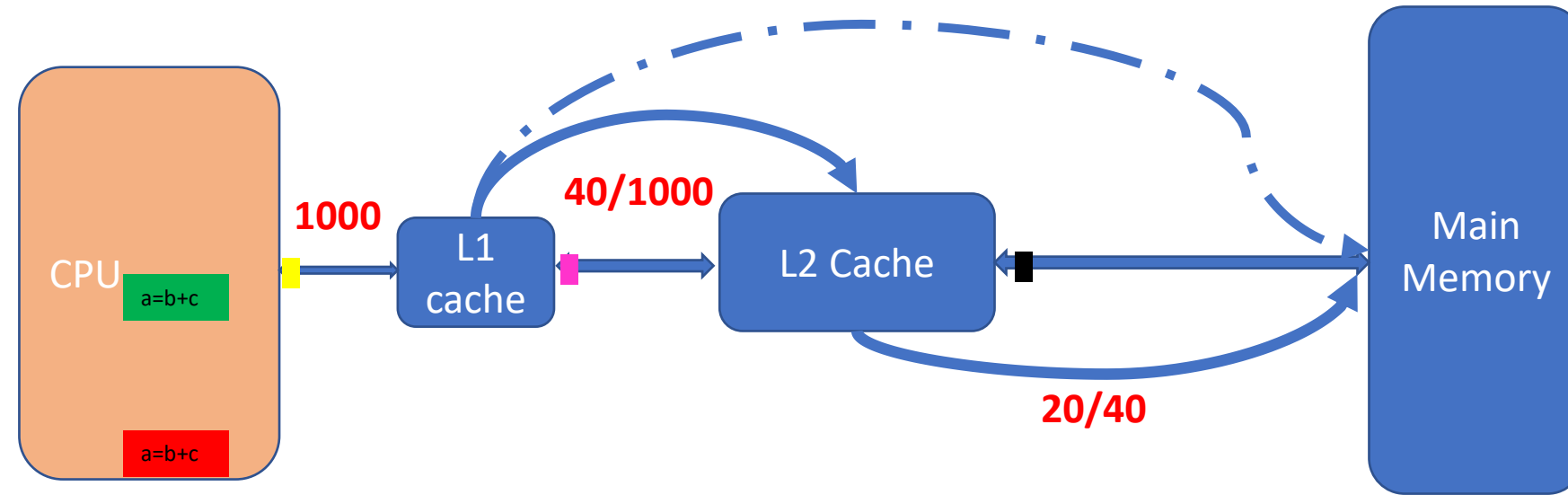
- Hit Time_{L2} = Item found in the Level 2 Cache
- Miss Rate_{L2} = How frequently element not found in Level 2 Cache
- Miss Penalty_{L2} = Extra time spent to bring item from RAM.

AMAT? For two Level Cache

$$\text{AMAT} = \text{Hit Time}_{L1} + \text{Miss Rate}_{L1} \times [\text{Hit Time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2}]$$

Microprocessor & Computer Architecture (μpCA)

4th Optimization: Multilevel Caches to Reduce Miss Penalty



Suppose that in 1000 memory references there are 40 misses in the first level cache and 20 misses in the second –level cache.

What is Miss Rate of Level 1 Cache? **40/1000**

What is Miss Rate of Level 2 Cache? **20/1000** or **20/40**

Microprocessor & Computer Architecture (μpCA)

4th Optimization: Multilevel Caches to Reduce Miss Penalty

Local miss rate:

- The number of misses in the cache divided by the total number of memory accesses to this cache.

Ex: For first level cache it is, **Miss Rate_{L1} (40/1000)**

For second level cache it is, **Miss Rate_{L2} (20/40)**

Global miss rate:

- The number of misses in the cache divided by the total number of memory accesses generated by the processor.

Ex: Global miss rate for level1 cache is still **Miss Rate_{L1} (40/1000)**

but, for level2 cache it is : **Miss Rate_{L1} x Miss Rate_{L2}**

$$(40/1000) \times (20/40) = (20/1000)$$

Microprocessor & Computer Architecture (μpCA)

4th Optimization: Multilevel Caches to Reduce Miss Penalty



Suppose that in 1000 memory references there are 40 misses in the first level cache and 20 misses in the second –level cache. ***What are the various miss rates?***

Assume the miss penalty from the L2 cache to memory is 200 clock cycles, the hit time of the L2 cache is 10 clock cycles, hit time for L1 cache is 1 clock cycle.

What is the average memory access time ?

Solution:

The Miss rate [global & local] for the 1st level cache is $40/1000 = 4\%$

The local miss rate for 2nd Level cache is $20/40 = 50\%$

The Global miss rate of 2nd Level cache is $20/1000 = 2\%$

$$\text{AMAT} = \text{Hit Time}_{L1} + \text{Miss Rate}_{L1} \times [\text{Hit Time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2}]$$

$$\text{AMAT} = 1 + 4\% \times [10 + 50\% \times 200]$$

$$= 1 + 0.04 \times [10 + 0.5 \times 200]$$

$$= 1 + 0.04 \times 110$$

$$= 5.4$$

Average Memory Stalls Per instruction = $(AMAT - \text{Hit time}_{L1}) \times \text{Average \# of memory references per instruction}$

or

Average memory stalls per instruction = $\text{Misses per instruction}_{L1} \times \text{Hit time}_{L2} + \text{Misses per instruction}_{L2} \times \text{Miss Penalty}_{L2}$

or

Average memory stalls per instruction = $(\text{Miss Rate}_{L1} \times \text{Hit time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2}) \times \text{Memory Reference per Instruction}$

or

Average memory stalls per instruction = $(\text{Miss Rate}_{L1} \times \text{Miss Penalty}_{L1} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2}) \times \text{Memory Reference per Instruction}$

Microprocessor & Computer Architecture (μpCA)

Average Memory Stalls Per Instruction



Suppose that in 1000 memory references there are 40 misses in the first level cache and 20 misses in the second –level cache. ***What are the various miss rates?***

Assume the miss penalty from the L2 cache to memory is 200 clock cycles, the hit time of the L2 cache is 10 clock cycles, hit time for L1 cache is 1 clock cycle and **there are 1.5 memory references per instruction.**

What is the Average Stall Cycles Per Instruction?

Solution:

$$\text{\# of misses per instruction} = \frac{\text{\# of memory references}}{\text{\# of memory references per instructions}} = \frac{1000}{1.5} = 667 \text{ ins.}$$

Thus, for L1 cache ,

**\# misses for 40 memory accesses for 1000 instructions = 40 x 1.5 = 60 misses and
for 20 misses for L2 cache it is 1.5 x 20 = 30 misses.**

Microprocessor & Computer Architecture (μpCA)

Average Memory Stalls Per Instruction



$$\begin{aligned}\text{Average memory stalls per instruction} &= \text{Misses per instruction}_{L1} \times \text{Hit time}_{L2} + \text{Misses per instruction}_{L2} \times \\ &\text{Miss Penalty}_{L2} \\ &= (60/1000) \times 10 + (30/1000) \times 200 \\ &= 0.06 \times 10 + 0.03 \times 200 = 0.06 + 6 \\ &= 6.6 \text{ clock cycles.}\end{aligned}$$

or

$$\begin{aligned}\text{Average memory stalls per instruction} &= (\text{AMAT} - \text{Hit time}_{L1}) \times \text{Memory references per instruction} \\ &= (5.4 - 1.0) \times 1.5 = 6.6 \text{ clock cycles.}\end{aligned}$$

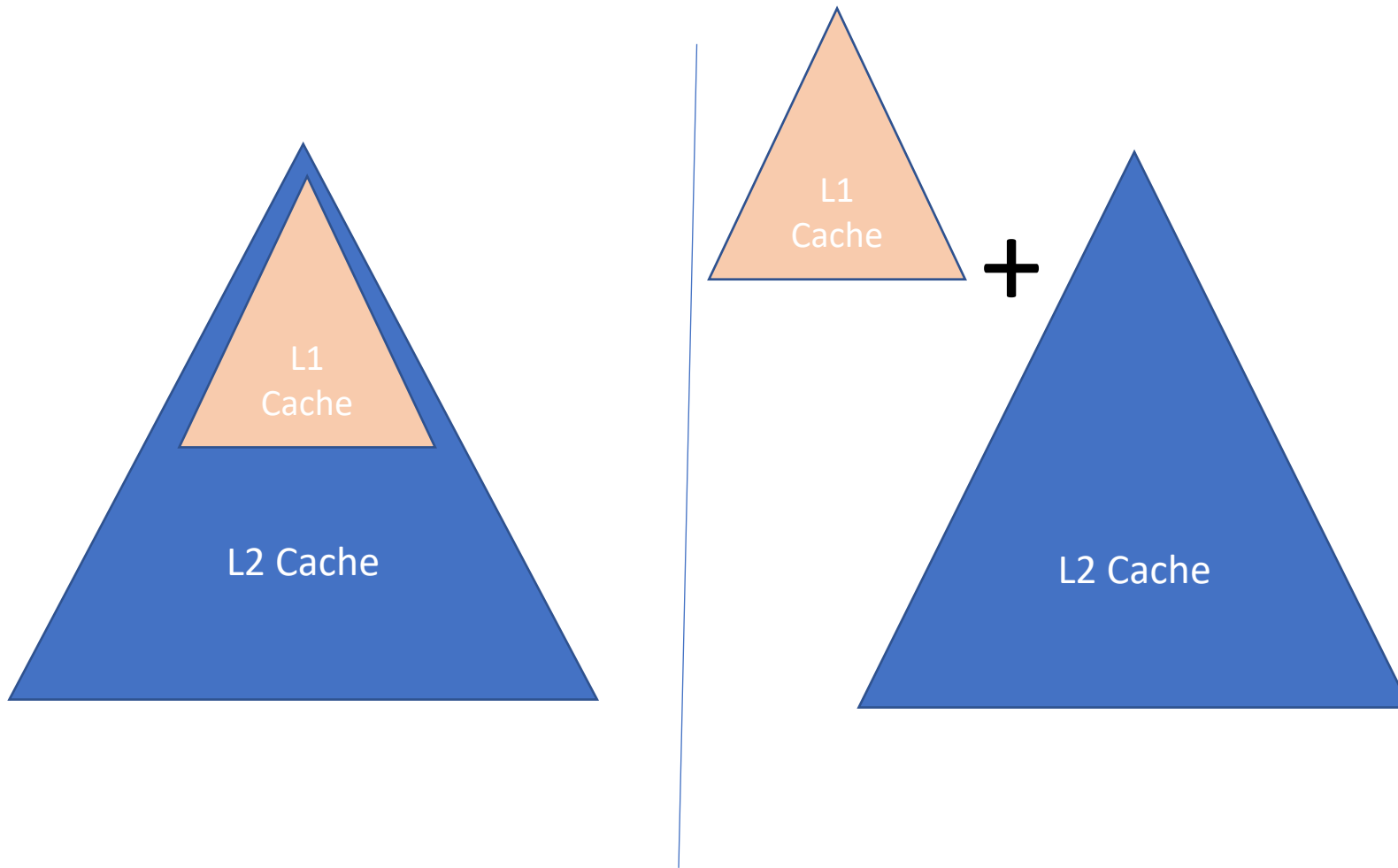
or

$$\begin{aligned}\text{Average memory stalls per instruction} &= (\text{Miss Rate}_{L1} \times \text{Hit time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2}) \times \text{Memory Reference per Instruction} \\ &= ((40/1000 \times 10) + (20/1000 \times 200)) \times 1.5 \\ &= 6.6 \text{ Clock Cycles}\end{aligned}$$

Microprocessor & Computer Architecture (μpCA)

Why Global Miss Rate of 2nd Level Cache is Important?

Which one is Correct?



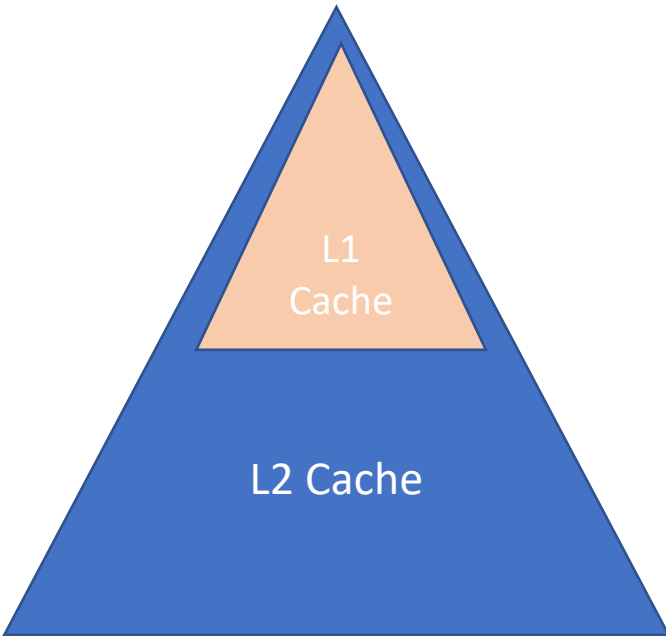
Why Global Miss Rate of 2nd Level Cache is Important?

First perspective :

Global cache miss rate is very similar to the single cache miss rate of the second level cache.

Second Perspective:

- Local cache miss rate is not a good measure of secondary caches.
- It is a function of the miss rate of the first level cache.
- Can vary by changing the first – level cache.



Note: Global cache miss rate should be used when evaluating second level caches.

Parameters for choosing 2nd Level Cache

- Miss Rate of 2nd level cache is function of 1st Level Cache.
- Speed of 2nd Level Cache will affect the Miss Penalty of 1st Level Cache.
- Design of 2nd Level cache should compensate all the deficiency in designing of 1st Cache.
- Thus, the strategies or Policy used in 2nd Level cache need not be same as design of 1st Level Cache.

Parameters for choosing 2nd Level Cache

Example Given the data below, what is the impact of second-level cache associativity on its miss penalty?

- Hit time_{L2} for direct mapped = 10 clock cycles.
- Two-way set associativity increases hit time by 0.1 clock cycle to 10.1 clock cycles.
- Local miss rate_{L2} for direct mapped = 25%.
- Local miss rate_{L2} for two-way set associative = 20%.

Answer For a direct-mapped second-level cache, the first-level cache miss penalty is

$$\text{Miss penalty}_{1\text{-way L2}} = 10 + 25\% \times 200 = 60.0 \text{ clock cycles}$$

Adding the cost of associativity increases the hit cost only 0.1 clock cycle, making the new first-level cache miss penalty:

$$\text{Miss penalty}_{2\text{-way L2}} = 10.1 + 20\% \times 200 = 50.1 \text{ clock cycles}$$

Microprocessor & Computer Architecture (μpCA)

5th Optimization: Giving priority to Read Misses over Write

Case 1: Write Back Policy

Scenario

- **P is referred by the Processor.**
- **P Should replace A.**
- **A is Dirty**

Known Solution

- **Write A back in Memory**
- **Replace P by A**

Cache

Valid Bit	Dirty Bit	TAG	Data
1	1	10100	A=11
1	0	10100	B=5
1	0	10100	C=6

Main Memory

10100.....1010..0000	A= 0
10100.....1010..0001	B=5
10100.....1010..0010	C=6
00010.....1010..0000	P=100

Optimized Solution (Give Priority for CPU Reference / Read)

- **Place A on to buffer to save time.**
- **Replace P by A, to provide quick access .**
- **Write A into memory parallelly when Processor is using P.**

Microprocessor & Computer Architecture (μpCA)

5th Optimization: Giving priority to Read Misses over Write

Case 2: Write Through

Scenario

- **P** is referred by the Processor.
- **P** Should replace **A**.

Cache

Valid Bit	TAG	Data
1	10100	A=11
1	10100	B=5
1	10100	C=6

		A=11
--	--	------

Write Buffer

Main Memory

10100.....1010..0000	A= 0
10100.....1010..0001	B=5
10100.....1010..0010	C=6
00010.....1010..0000	P=100

Microprocessor & Computer Architecture (μpCA)

5th Optimization: Giving priority to Read Misses over Write

Case 2: Write Through

Scenario

- **P is referred by the Processor.**
- **P Should replace A.**

Known Solution

- **Place A on the Write Buffer**
- **Replace P by A**

What may go wrong?

If A is referred again!

Cache

Valid Bit	TAG	Data
1	00010	P=100
1	10100	B=5
1	10100	C=6

		A=11
--	--	------

Write Buffer

Main Memory

10100.....1010..0000	A= 0
10100.....1010..0001	B=5
10100.....1010..0010	C=6
00010.....1010..0000	P=100

Solution1: Wait till A is written back in Memory

Solution2: (Give Priority for CPU Reference / Read)

Search for A in buffer, if found compare A in Memory,

- i. if both are same, Replace.
- ii. Otherwise, read from Buffer

Consider the following code sequence.

Ex: STR R3, [R0,512]

LDR R1, [R0,1024]

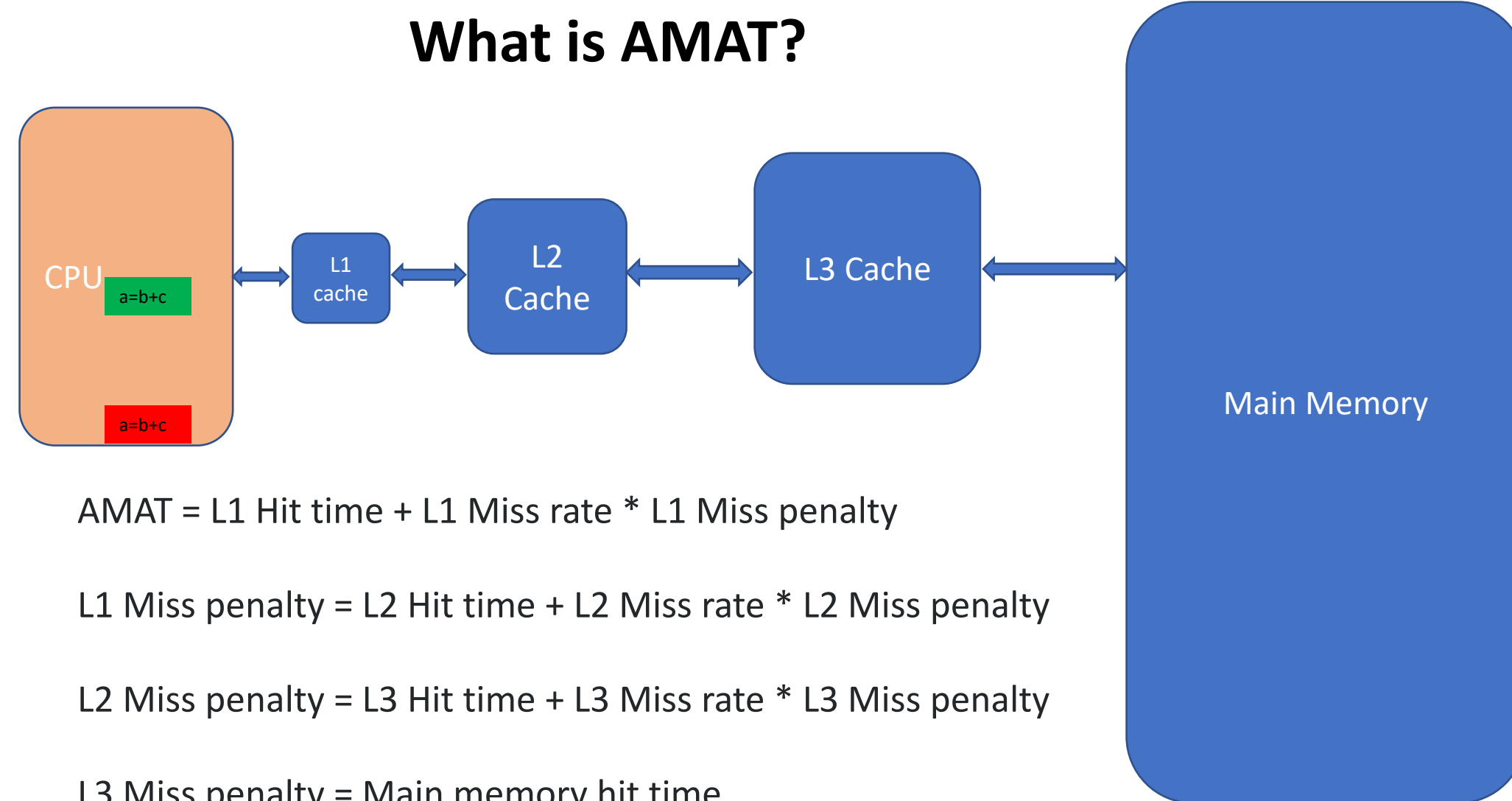
LDR R2, [R0,512] Assume Direct mapped:

Write-through cache that maps 512 and 1024 to the same block. Four word write buffer that is not checked on a read miss. Will the value in R2 always be equal to the value in R3? o R2!

Ans:

- This is a read-after-write data hazard in memory.
- The data in R3 are placed into the write buffer after the STR.
- The following LDR instruction uses the same cache index and is therefore a miss.
- The second LDR instruction, tries to put the value in location 512 into the register R2.
- This also results in a miss.
- If the write buffer hasn't completed writing to location 512 in memory,
- The read of location 512 will put the old, wrong value into the cache block and then into R2.
- Without proper precautions, R3 would not be equal to R2!

What is AMAT?



Improving Hit Time



THANK YOU

Dr. D. C. Kiran

Department of Computer Science and Engineering

dckiran@pes.edu

9829935135