

# CS223 : Computer Architecture & Organization

**Lecture 23 [01.04.2022]**

## **Advanced Cache Optimizations-Tutorials**



**Dr. John Jose**

**Associate Professor**

**Department of Computer Science & Engineering  
Indian Institute of Technology Guwahati, Assam.**

# Multilevel Caches

Assume a 2-level cache system with the following specifications. L1 Hit Time = 1 cycle, L1 Miss Rate = 2.5%, L2 Hit Time = 6 cycles, L2 Miss Rate = 17% (% L1 misses that miss), L2 Miss Penalty = 120 cycles. Compute the average memory access time.

$$\begin{aligned} \text{AMAT} &= \text{ht}_1 + \text{mr}_1 \times \text{MP} \\ &= \text{ht}_1 + \text{mr}_1 \times (\text{ht}_2 + \text{mr}_2 \times \text{MP}) \\ &= 1 + 0.025 \times (6 + 0.17 \times 120) = 1.66 \text{ CC} \end{aligned}$$

# Optimization

A cache has access time (hit latency) of 10 ns and miss rate of 5%. An optimization was made to reduce the miss rate to 3% but the hit latency was increased to 15 ns. Under what condition this change will result in better performance (Lower AMAT)?

# Optimization

A cache has access time (hit latency) of 10 ns and miss rate of 5%. An optimization was made to reduce the miss rate to 3% but the hit latency was increased to 15 ns. Under what condition this change will result in better performance (Lower AMAT)?

$$\text{AMAT}_1 = \text{HT}_1 + \text{MR}_1 \times \text{MP}$$

$$\text{HT}_1 = 10\text{ns}; \text{MR}_1 = 0.05$$

$$\text{AMAT}_2 = \text{HT}_2 + \text{MR}_2 \times \text{MP}$$

$$\text{HT}_2 = 15\text{ns}; \text{MR}_2 = 0.03$$

$$\text{AMAT}_2 < \text{AMAT}_1$$

$$15 + 0.03 \times \text{MP} < 10 + 0.05 \times \text{MP}$$

$$5 < 0.02\text{MP} \rightarrow \text{MP} > 250 \text{ ns}$$

# Optimization

A cache has hit rate of 95%, block size of 128B, cache hit latency of 5ns. Main memory takes 50 ns to return first word (32 bits) of a block and 10 ns for each subsequent word.

- (a) What is the miss latency of the cache?
- (b) If doubling the cache block size reduces the miss rate to 3%, does it reduce AMAT?

# Optimization

hit rate of 95%, 28B blocks, cache hit latency of 5ns. Main memory takes 50 ns to return first word (32 bits) of a block and 10 ns for each subsequent word. (a) What is the miss latency of the cache?

(b) If doubling the cache block size reduces the miss rate to 3%, does it reduces AMAT?

$H_r = 0.95$ ;  $BS = 128B$ ;  $H_t = 5 \text{ ns}$ ;  $1 \text{ word} = 4B \text{ (32 bits)}$

$\# \text{ words/ block} = 128B / 4B = 32$

(A)  $MP = 50 + (31 \times 10) = 360 \text{ ns}$

$AMAT_1 = 5 + 0.05 \times 360 = 23 \text{ ns}$

(B)  $\# \text{ words/ block} = 256B / 4B = 64$  ;

$MP = 50 + 63 \times 10 = 680 \text{ ns}$

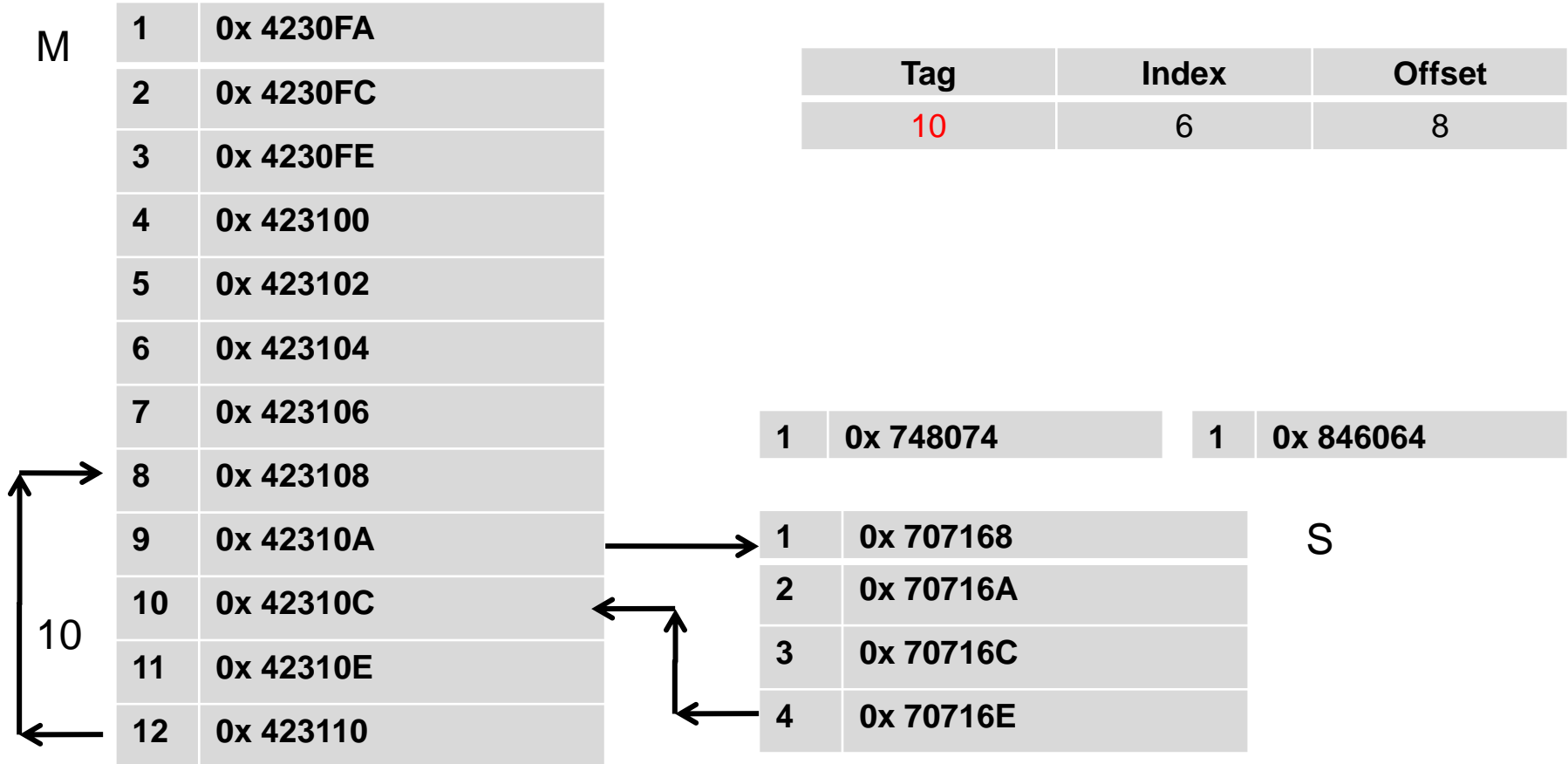
$AMAT_2 = 5 + 0.03 \times 680 = 25.4 \text{ ns}$

Doubling block size will not reduce AMAT

# Cache Mapping

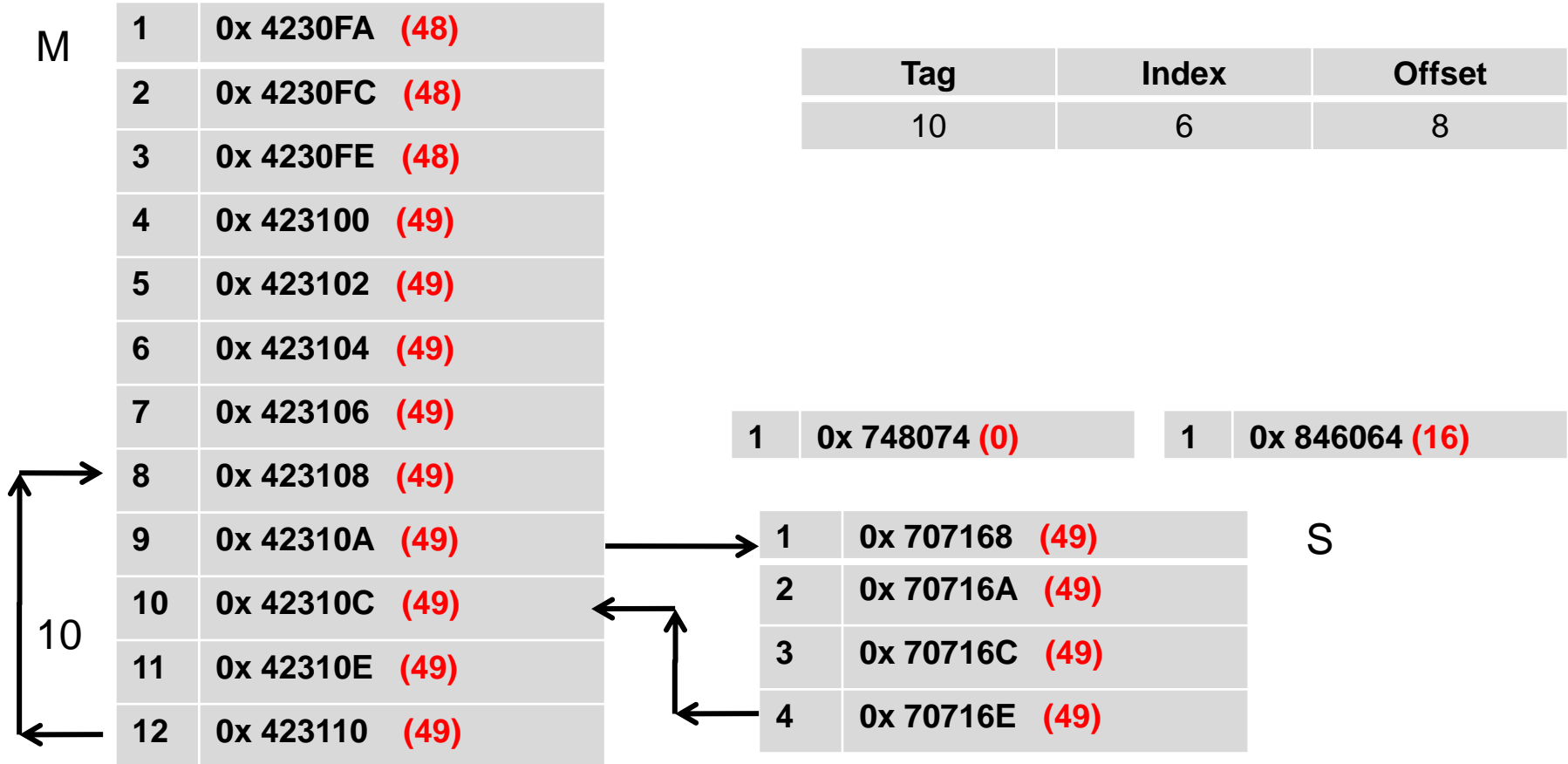
A 16KB direct mapped 256B block unified cache is attached to a 16MB main memory system. The word length as well as instruction length of the processor is 16 bits. Consider a program that consists of a main routine M which in turn calls a subroutine S. M consists of 12 instruction words which are loaded in the main memory from the address 0x4230FA onwards. The last five instructions of M is a loop that is iterated 10 times. The second instruction in the loop is a call to subroutine S. S consists of 4 instruction words loaded in the main memory from the address 0x70F168. The last instruction of S is a subroutine return back to M. The only two data words that are used by M and S are at addresses 0x748074 and 0x846064. Assume the caches are initially empty. Ignore OS level interruption and subsequent cache impact on context switching. List out the block numbers (in decimal) in the cache that are non-empty after the execution of the program.

# Cache Mapping by OS





# Cache Mapping by OS



# Cache Mapping by OS

Find the number of cache misses occurred during the execution of the program.

M1

M4

S1, M10, S1, M10, ..... (10 TIMES) = 22 MISSES

How many cache block evictions happened during the execution of the program?

20 EVICTIONS

List out the block numbers (in decimal) in the cache that are non-empty after the execution of the program.

ALL BLOCKS EXCEPT 0,16,48,49

# Reference

- ❖ **Computer Architecture-A Quantitative Approach** (5th edition),  
John L. Hennessy, David A. Patterson, Morgan Kaufman.
- ❖ Chapter 2: Memory Hierarchy Design
  - ❖ Section 2.1: Introduction
  - ❖ Section 2.2: Ten Advanced Cache Optimizations of Cache Performance
- ❖ **NPTEL Video Links:**
  - ❖ <https://tinyurl.com/yf94dnby>
  - ❖ <https://tinyurl.com/yfgmfsmm>



**johnjose@iitg.ac.in**  
**<http://www.iitg.ac.in/johnjose/>**