

# CS510 Computer Architecture

## Lecture 14: Cache Write Policies & 2-Level Caches

**Soontae Kim**

**Spring 2017**

**School of Computing, KAIST**

# Notice

- **Term project proposal**
  - On May 10 (Wednesday)
  - Prepare less than 10 slides in 5 minutes.
    - **No proposal report required**
  - All team members must prepare parts of presentation.
  - Any topic related to computer architecture; if you are not sure that your topic is appropriate, you can discuss with me.
  - Practice your presentation before coming to class several times so that you can finish your presentation in time; we have only 75 minutes for all of your presentations.
  - You may need to change your topic or direction if your topic conflicts with other team's topic or you find a difficulty in performing your projects.

# Cache Write Miss Policy

- Since data are usually not needed immediately on a write miss, two options exist on a cache write miss:

## Write Allocate:

The cache block is loaded on a write miss followed by write hit actions.

## No-Write Allocate:

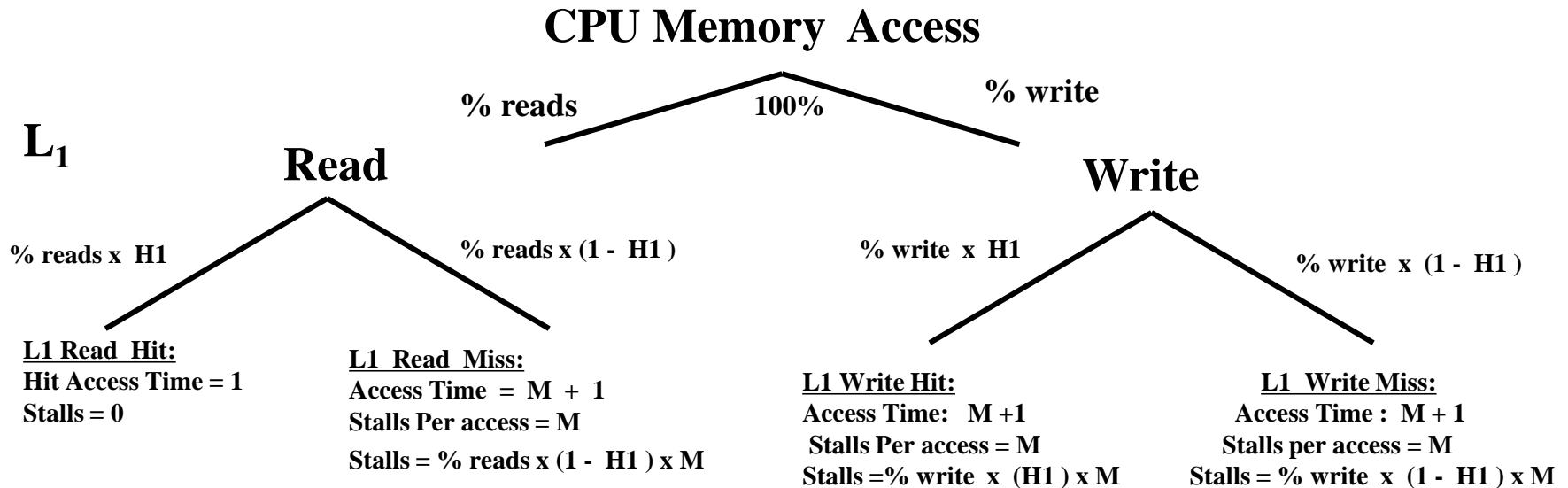
The block is modified in the lower level (lower cache level, or main memory) and not loaded into cache.

*While any of the above two write miss policies can be used with either write back or write through:*

- Write back caches always use write allocate to capture subsequent writes to the block in cache.
- Write through caches usually use no-write allocate since subsequent writes still have to go to memory.

# Memory Access Tree, Unified L<sub>1</sub>

**Write Through, No Write Allocate, No Write Buffer**



$$\text{Stall Cycles Per Memory Access} = \% \text{ reads} \times (1 - H1) \times M + \% \text{ write} \times M$$

$$\text{AMAT} = 1 + \% \text{ reads} \times (1 - H1) \times M + \% \text{ write} \times M$$

$$\text{CPI} = \text{CPI}_{\text{execution}} + (1 + \text{fraction of loads/stores}) \times \text{Stall Cycles per access}$$

$M$  = Miss Penalty  
 $H1$  = Level 1 Hit Rate  
 $1 - H1$  = Level 1 Miss Rate

# Reducing Write Stalls For Write Through Cache using Write Buffers

- To reduce write stalls when write through is used, a write buffer is used to eliminate or reduce write stalls:

- Perfect write buffer: All writes are handled by write buffer, no stalling for writes

- In this case (for unified L1 cache):

$$\text{Stall Cycles Per Memory Access} = \% \text{ reads} \times (1 - H1) \times M$$

(i.e No stalls at all for writes)

- Realistic Write buffer: Write stalls are not eliminated when the write buffer is full.

- In this case (for unified L1 cache):

$$\text{Stall Cycles/Memory Access} = ( \% \text{ reads} \times (1 - H1) + (\% \text{ write stalls not eliminated}) ) \times M$$

# Write Through Cache Performance Example

- A CPU with  $CPI_{\text{execution}} = 1.1$  Mem accesses per instruction = 1.3
- Uses a unified L1 Write Through, No Write Allocate, with:
  - No write buffer.
  - Perfect Write buffer
  - A realistic write buffer that eliminates 85% of write stalls
- Instruction mix: 50% arith/logic, 15% load, 15% store, 20% control
- Assume a cache miss rate of 1.5% and a miss penalty of 50 cycles.

$$CPI = CPI_{\text{execution}} + \text{mem stalls per instruction}$$

$$\% \text{ reads} = 1.15/1.3 = 88.5\% \quad \% \text{ writes} = .15/1.3 = 11.5\%$$

## With No Write Buffer :

$$\text{Mem Stalls/ instruction} = 1.3 \times 50 \times (88.5\% \times 1.5\% + 11.5\%) = 8.33 \text{ cycles}$$

$$CPI = 1.1 + 8.33 = 9.43$$

## With Perfect Write Buffer (all write stalls eliminated):

$$\text{Mem Stalls/ instruction} = 1.3 \times 50 \times (88.5\% \times 1.5\%) = 0.86 \text{ cycles}$$

$$CPI = 1.1 + 0.86 = 1.96$$

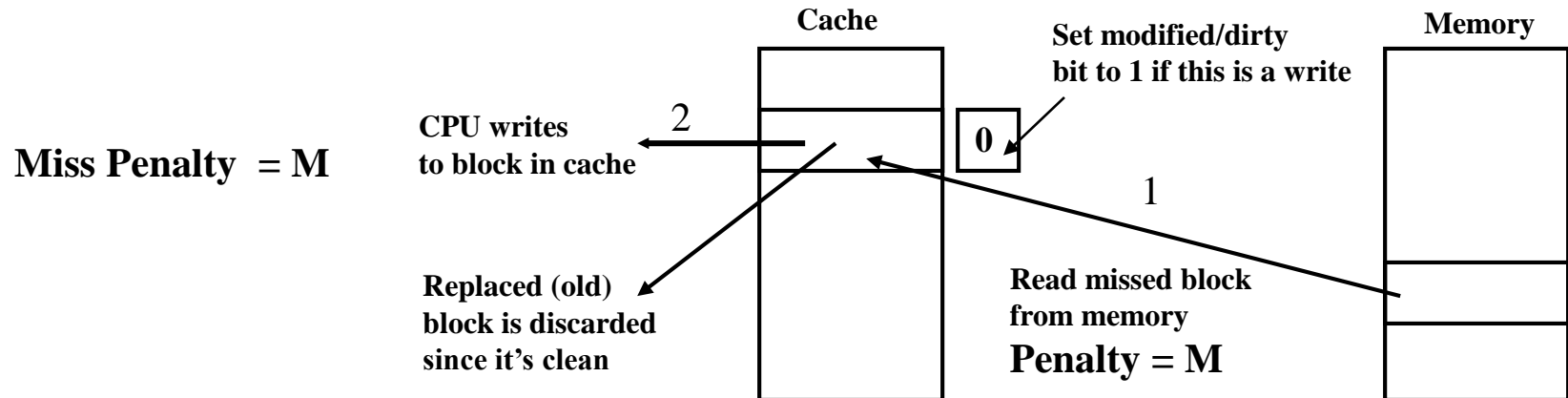
## With Realistic Write Buffer (eliminates 85% of write stalls)

$$\text{Mem Stalls/ instruction} = 1.3 \times 50 \times (88.5\% \times 1.5\% + 15\% \times 11.5\%) = 1.98 \text{ cycles}$$

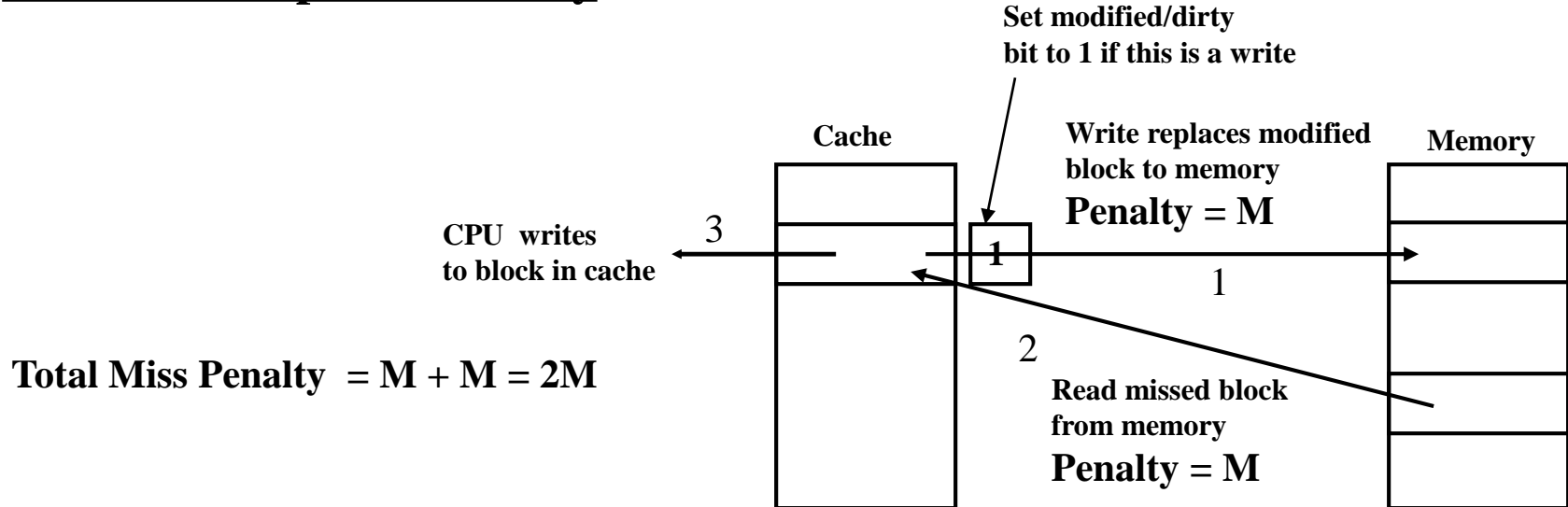
$$CPI = 1.1 + 1.98 = 3.08$$

# Write Back Cache With Write Allocate: Cache Miss Operation

## Block to be replaced is clean



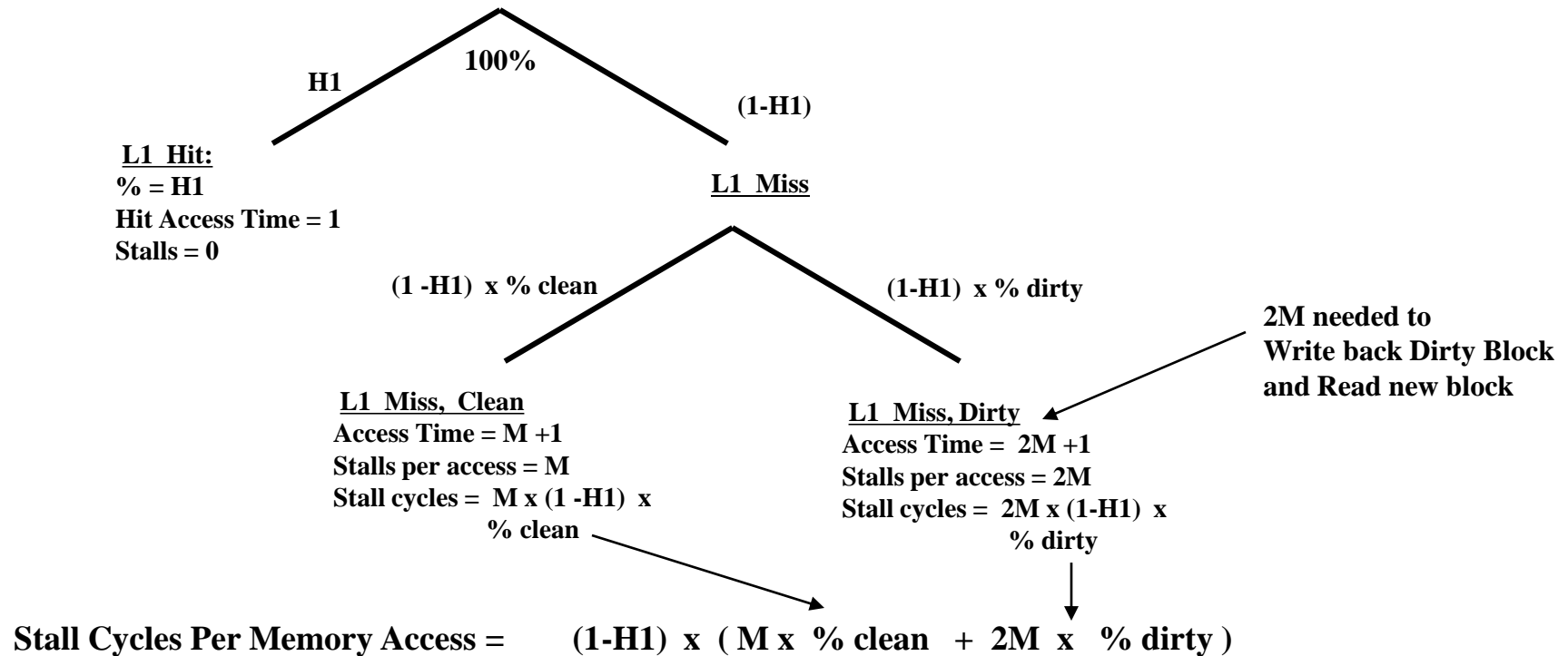
## Block to be replaced is dirty



**M = Miss Penalty = stall cycles per access resulting from missing in cache**

# Memory Access Tree Unified L<sub>1</sub> Write Back, With Write Allocate

## CPU Memory Access



$$\text{AMAT} = 1 + \text{Stall Cycles Per Memory Access}$$

$$\text{CPI} = \text{CPI}_{\text{execution}} + (1 + \text{fraction of loads/stores}) \times \text{Stall Cycles per access}$$



# Write Back Cache Performance Example

- A CPU with  $CPI_{\text{execution}} = 1.1$  uses a unified L1 with write back, with write allocate, and the probability a cache block is dirty = 10%
- Instruction mix: 50% arith/logic, 15% load, 15% store, 20% control
- Assume a cache miss rate of 1.5% and a miss penalty of 50 cycles.

$$CPI = CPI_{\text{execution}} + \text{mem stalls per instruction}$$

$$\text{Mem Stalls per instruction} =$$

$$\text{Mem accesses per instruction} \times \text{Stalls per access}$$

$$\text{Mem accesses per instruction} = 1 + .3 = 1.3$$

$$\text{Stalls per access} = (1-H1) \times (M \times \% \text{ clean} + 2M \times \% \text{ dirty})$$

$$\text{Stalls per access} = 1.5\% \times (50 \times 90\% + 50 \times 2 \times 10\%) = .825 \text{ cycles}$$

$$\text{Mem Stalls per instruction} = 1.3 \times .825 = 1.07 \text{ cycles}$$

$$AMAT = 1 + 0.825 = 1.825 \text{ cycles}$$

$$CPI = 1.1 + 1.07 = 2.17$$

The ideal CPU with no misses is  $2.17/1.1 = 1.97$  times faster

# 2-Level Cache: $L_1, L_2$

CPU

$L_1$  Cache

Hit Rate =  $H_1$ ,  
Hit Access Time = 1 cycle (No Stall)  
Stalls for hit access =  $T_1 = 0$

$L_2$  Cache

Local Hit Rate =  $H_2$   
Stalls per hit access =  $T_2$   
Hit Access Time =  $T_2 + 1$  cycles

Main Memory

Memory access penalty,  $M$   
(stalls per main memory access)  
Access Time =  $M + 1$

## Goal of multi-level Caches:

Reduce the effective miss penalty incurred by level 1 cache misses using additional levels of cache that capture some of these misses

# Miss Rates For Multi-Level Caches

- **Local Miss Rate:** This rate is the number of misses in a cache level divided by the number of memory accesses to this level. **Local Hit Rate = 1 - Local Miss Rate**
- **Global Miss Rate:** The number of misses in a cache level divided by the total number of memory accesses generated by the CPU.
- **Since level 1 receives all CPU memory accesses, for level 1:**

$$\text{Local Miss Rate} = \text{Global Miss Rate} = 1 - H_1$$

- **For level 2 since it only receives those accesses missed in level 1:**

$$\text{Local Miss Rate} = \text{Miss rate}_{L_2} = 1 - H_2$$

$$\begin{aligned}\text{Global Miss Rate} &= \text{Miss rate}_{L_1} \times \text{Miss rate}_{L_2} \\ &= (1 - H_1) \times (1 - H_2)\end{aligned}$$

## 2-Level Cache Performance (Ignoring Write Policy)

$$\text{CPUtime} = \text{IC} \times (\text{CPI}_{\text{execution}} + \text{Mem Stall cycles per instruction}) \times \text{C}$$

$$\text{Mem Stall cycles per instruction} = \text{Mem accesses per instruction} \times \text{Stall cycles per access}$$

- For a system with 2 levels of unified cache, assuming no penalty when found in L<sub>1</sub> cache:

Stall cycles per memory access =

$$[\text{miss rate } L_1] \times [\text{Hit rate } L_2 \times \text{Hit time } L_2 + \text{Miss rate } L_2 \times \text{Memory access penalty}] =$$

$$(1-H_1) \times H_2 \times T_2 + (1-H_1)(1-H_2) \times M$$

L1 Miss, L2 Hit

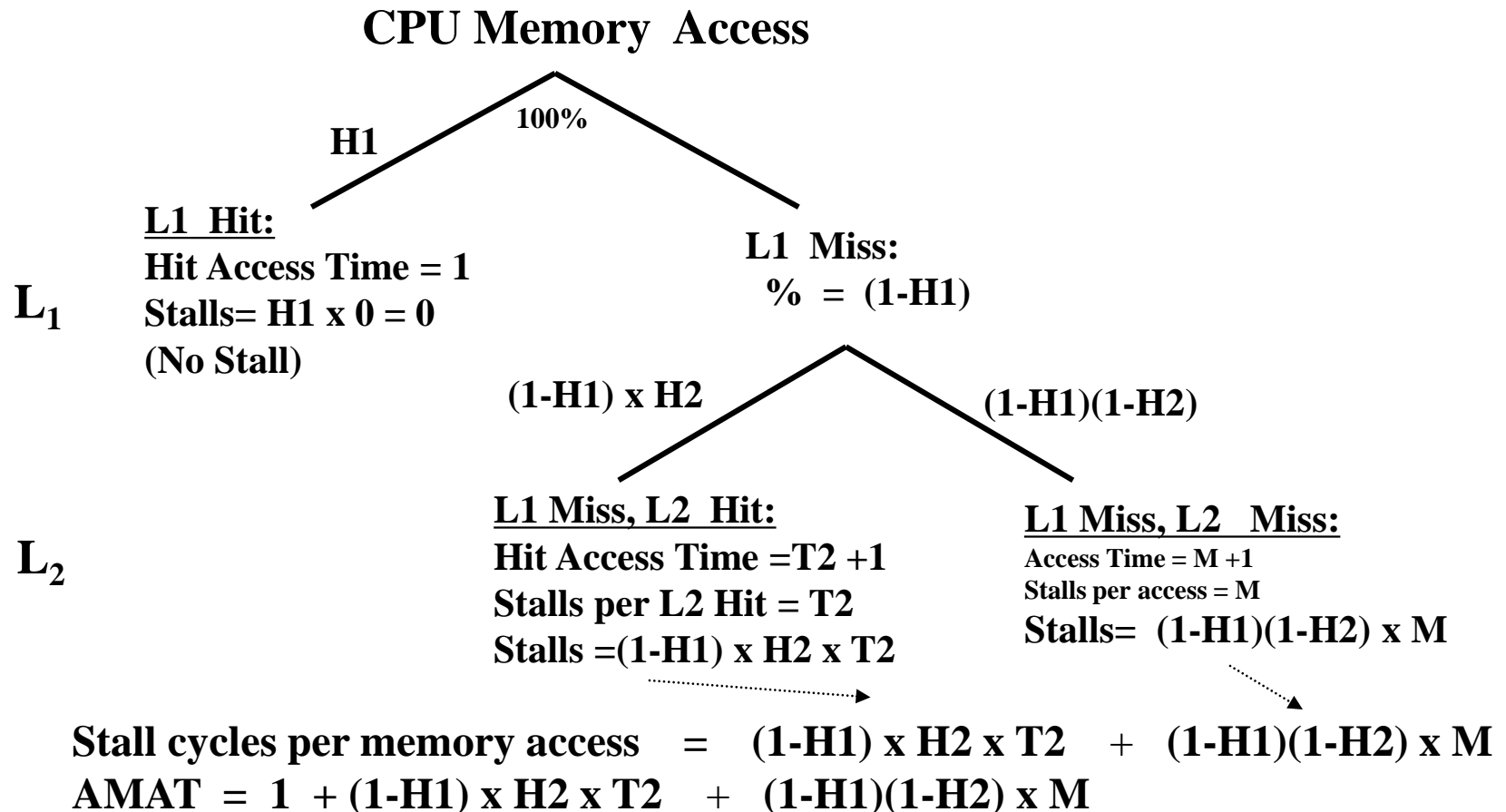
H<sub>1</sub> = L1 Hit Rate  
H<sub>2</sub> = Local L2 Hit Rate  
T<sub>2</sub> = stall cycles per L2 hit

L1 Miss, L2 Miss:  
Must Access Main Memory

# 2-Level Cache (Both Unified) Performance

## Memory Access Tree (Ignoring Write Policy)

### CPU Stall Cycles Per Memory Access



# Two-Level Cache Example

- CPU with  $CPI_{\text{execution}} = 1.1$  running at clock rate = 500 MHz
- 1.3 memory accesses per instruction.
- $L_1$  hit access time = 1 cycle (no stall on a hit), a miss rate of 5%
- $L_2$  hit access time = 3 cycles ( $T_2 = 2$  stall cycles per hit) with local miss rate 40%,
- Memory access penalty,  $M = 100$  cycles (stalls per access). Find CPI ...

$$CPI = CPI_{\text{execution}} + \text{Mem Stall cycles per instruction}$$

$$\text{With No Cache, } CPI = 1.1 + 1.3 \times 100 = 131.1$$

$$\text{With single } L_1, \quad CPI = 1.1 + 1.3 \times .05 \times 100 = 7.6$$

$$\begin{aligned} \text{Stall cycles per memory access} &= (1-H_1) \times H_2 \times T_2 + (1-H_1)(1-H_2) \times M \\ &= .05 \times .6 \times 2 + .05 \times .4 \times 100 \\ &= .06 + 2 = 2.06 \end{aligned}$$

$$AMAT = 2.06 + 1 = 3.06$$

$$\begin{aligned} \text{Mem Stall cycles per instruction} &= \text{Mem accesses per instruction} \times \text{Stall cycles per access} \\ &= 1.3 \times 2.06 = 2.678 \end{aligned}$$

$$CPI = 1.1 + 2.678 = 3.778$$

$$\text{Speedup} = 7.6/3.778 = 2$$

# Write Policy For 2-Level Cache

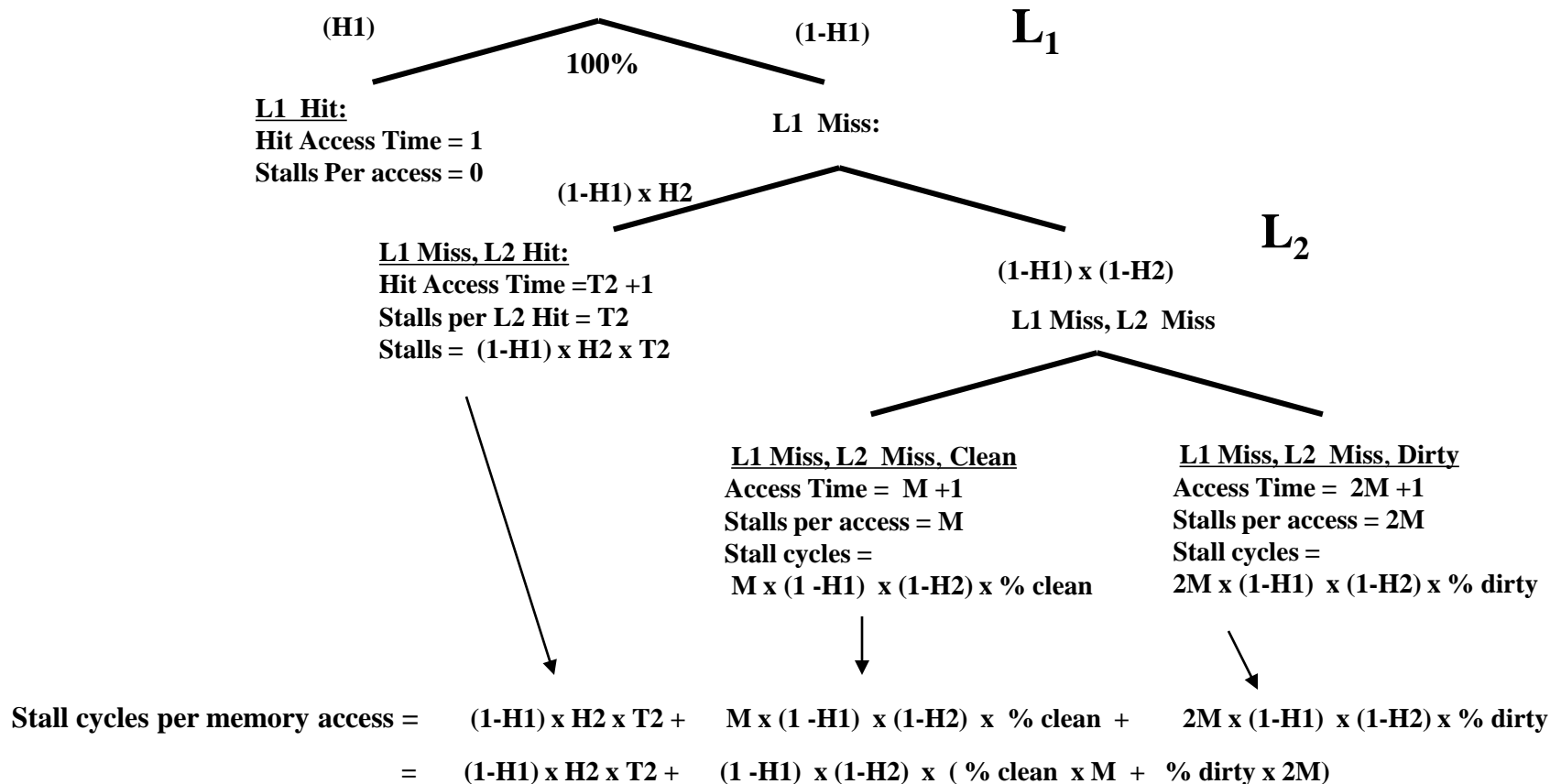
- Write Policy For Level 1 Cache:
  - Usually Write through to Level 2
  - Write allocate is used to reduce level 1 miss reads.
  - Use write buffer to reduce write stalls
- Write Policy For Level 2 Cache:
  - Usually write back with write allocate is used.
    - To minimize memory bandwidth usage.
- The above 2-level cache write policy results in **inclusive L2 cache** since the content of L1 is also in L2
  - Common in the majority of CPUs with 2-levels of cache
  - As opposed to **exclusive L1, L2** (e.g AMD Athlon XP)

# 2-Level (Both Unified) Memory Access Tree

**L1: Write Through to L2, Write Allocate, With Perfect Write Buffer**

**L2: Write Back with Write Allocate**

## CPU Memory Access





# Two-Level Cache Example With Write Policy

- CPU with  $CPI_{\text{execution}} = 1.1$  running at clock rate = 500 MHz
- 1.3 memory accesses per instruction.
- For  $L_1$  :
  - Cache operates at 500 MHz (no stall on L1 Hit) with a miss rate of  $1-H1 = 5\%$
  - Write through to  $L_2$  with perfect write buffer with write allocate
- For  $L_2$ :
  - Hit access time = 3 cycles ( $T2 = 2$  stall cycles per hit) local miss rate  $1-H2 = 40\%$
  - Write back to main memory with write allocate
  - Probability a cache block is dirty = 10%
- Memory access penalty,  $M = 100$  cycles. Find CPI.
- Stall cycles per memory access =  $(1-H1) \times H2 \times T2 +$   
 $(1-H1) \times (1-H2) \times (\% \text{ clean} \times M + \% \text{ dirty} \times 2M)$   
 $= .05 \times .6 \times 2 + .05 \times .4 \times (.9 \times 100 + .1 \times 200)$   
 $= .06 + 0.02 \times 110 = .06 + 2.2 = 2.26$

- $AMAT = 2.26 + 1 = 3.26$  cycles

Mem Stall cycles per instruction = Mem accesses per instruction  $\times$  Stall cycles per access

$$= 1.3 \times 2.26 = 2.938$$

$$CPI = 1.1 + 2.938 = 4.038$$