

# CS510 Computer Architecture

## Lecture 13: Review: Cache Memory II

**Soontae Kim**

**Spring 2017**

**School of Computing, KAIST**

# Notice

- **Term project proposal**
  - On May 9 (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.

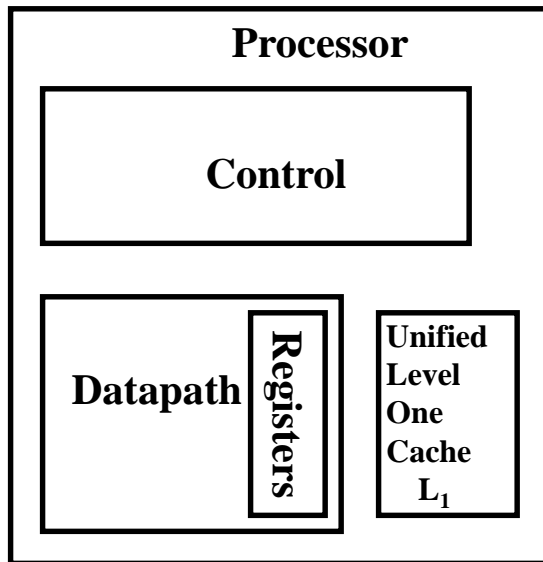
# Unified vs. Separate Level 1 Cache

- Unified Level 1 Cache (Princeton Memory Architecture).

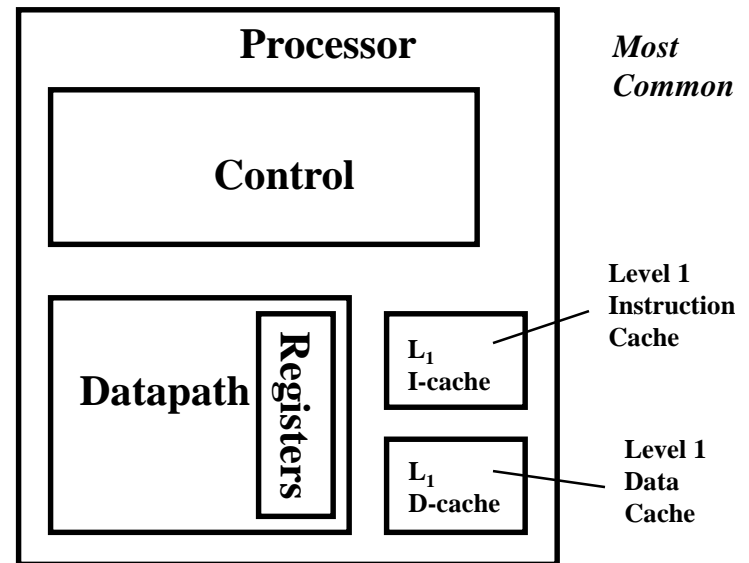
A single level 1 ( $L_1$ ) cache is used for both instructions and data.

- Separate instruction/data Level 1 caches (Harvard Memory Architecture):

The level 1 ( $L_1$ ) cache is split into two caches, one for instructions (instruction cache,  $L_1$  I-cache) and the other for data (data cache,  $L_1$  D-cache).



Unified Level 1 Cache  
(Princeton Memory Architecture)



Separate (Split) Level 1 Caches  
(Harvard Memory Architecture)

# ***Memory Hierarchy Performance:***

## **Average Memory Access Time (AMAT), Memory Stall cycles**

- **The Average Memory Access Time (AMAT):** The number of cycles required to complete an average memory access request by the CPU.
- **Memory stall cycles per memory access:** The number of stall cycles added to CPU execution cycles for one memory access.
- **Memory stall cycles per average memory access = (AMAT -1)**
- **For ideal memory: AMAT = 1 cycle, this results in zero memory stall cycles.**
- **Memory stall cycles per average instruction =**

$$\begin{array}{l} \text{Instruction} \\ \text{Fetch from \$} \end{array} \begin{array}{l} \nearrow \\ \rightarrow \end{array} \begin{array}{l} \text{Number of memory accesses per instruction} \\ \text{x Memory stall cycles per average memory access} \end{array} \\ = ( 1 + \text{fraction of loads/stores} ) \times ( \text{AMAT} - 1 )$$

$$\text{Base CPI} = \text{CPI}_{\text{execution}} = \text{CPI with ideal memory}$$

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

# Cache Performance:

## Single Level L1 Princeton (Unified) Memory Architecture

$$\text{CPUtime} = \text{Instruction count} \times \text{CPI} \times \text{Clock cycle time}$$

$$\text{CPI}_{\text{execution}} = \text{CPI with ideal memory}$$

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

$$\text{Mem Stall cycles per instruction} =$$

$$\text{Memory accesses per instruction} \times \text{Memory stall cycles per memory access}$$

Assuming no stall cycle on a cache hit (cache access time = 1 cycle, stall = 0)

$$\text{Cache Hit Rate} = H1 \quad \text{Miss Rate} = 1 - H1$$

$$\text{Memory stall cycles per memory access} = \text{Miss rate} \times \text{Miss penalty}$$

$$\text{AMAT} = 1 + \text{Miss rate} \times \text{Miss penalty}$$

$$\text{Memory accesses per instruction} = (1 + \text{fraction of loads/stores})$$

Miss Penalty = M = the number of stall cycles resulting from missing in cache

Thus for a unified L1 cache with no stall on a cache hit:

$$\text{CPI} = \text{CPI}_{\text{execution}} + (1 + \text{fraction of loads/stores}) \times (1 - H1) \times M$$

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

$$\text{AMAT} = \text{hit rate} \times \text{hit time} + (1 - \text{hit rate}) \times (\text{hit time} + \text{miss penalty})$$

$$= 1 + (1 - \text{hit rate}) \times \text{miss penalty}$$

$$\text{hit time} = \text{RAM access time} + \text{time to determine hit/miss}$$

$$\text{miss penalty} = \text{time to deliver missed block to the processor} + \text{time to replace a block in upper level cache}$$

$$\text{miss penalty} = \text{hit time} + \text{miss penalty} + \text{time to determine hit/miss} + \text{lower level cache} + \text{hit time}$$

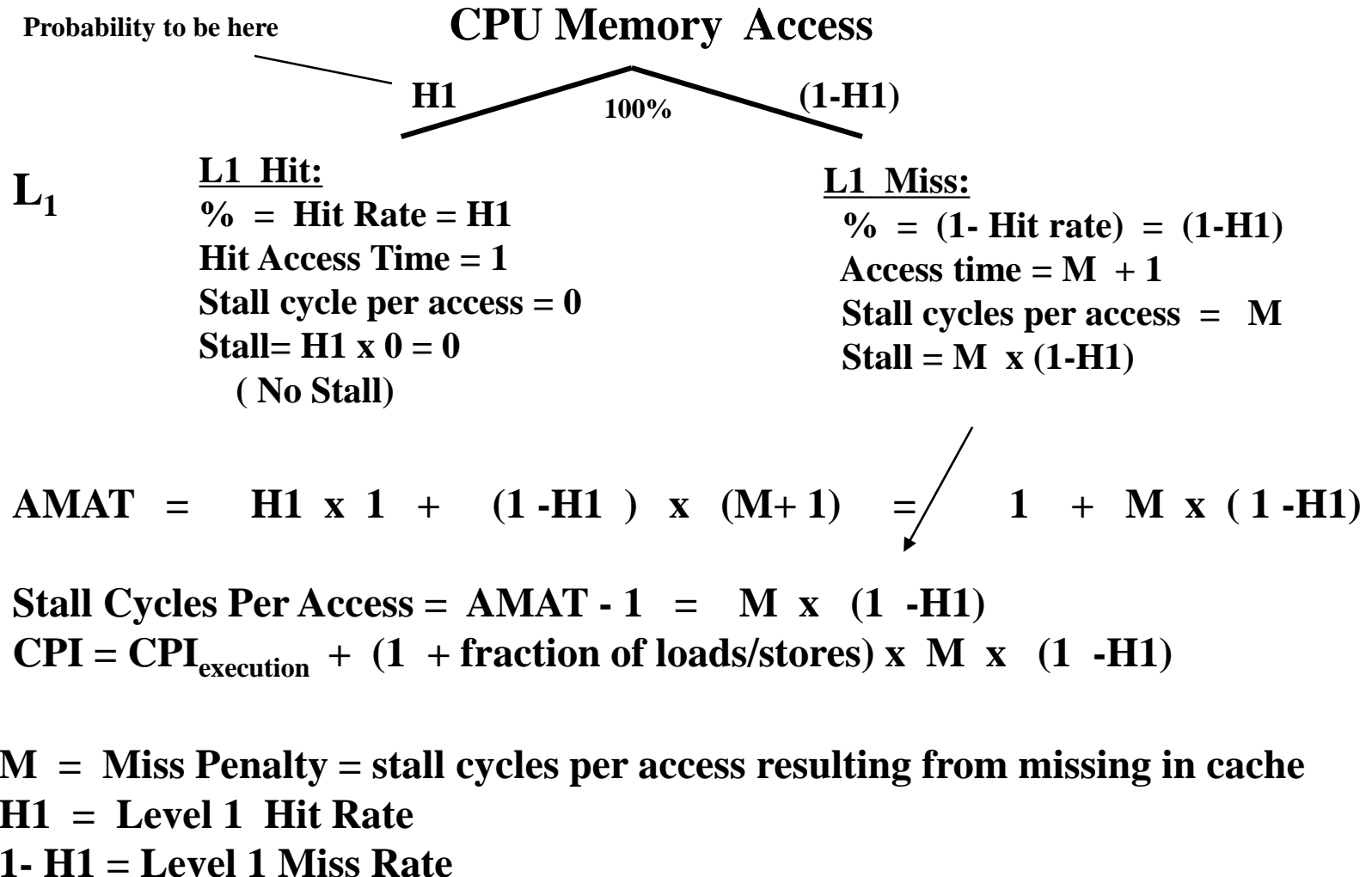
$$(\text{time to deliver missed block to the processor} + \text{time to replace a block in upper level cache} + \text{RAM access time})$$

$$\text{miss penalty} = \text{hit time} + \text{miss penalty} + \text{time to determine hit/miss} + \text{lower level cache} + \text{hit time}$$

cpu

# Memory Access Tree:

## For Unified Level 1 Cache



# Cache Performance Example

- Suppose a CPU executes at Clock Rate = 200 MHz (5 ns per cycle) with a single level of cache.
- $CPI_{\text{execution}} = 1.1$
- Instruction mix: 50% arith/logic, 30% load/store, 20% control
- Assume a cache miss rate of 1.5% and a miss penalty of  $M = 50$  cycles.

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

Mem Stalls per instruction =

$$\text{Mem accesses per instruction} \times \text{Miss rate} \times \text{Miss penalty}$$

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

Instruction fetch

Load/store

$$\text{Mem Stalls per memory access} = (1 - H1) \times M = .015 \times 50 = .75 \text{ cycles}$$

$$AMAT = 1 + .75 = 1.75 \text{ cycles}$$

$$\text{Mem Stalls per instruction} = 1.3 \times .015 \times 50 = 0.975$$

$$CPI = 1.1 + .975 = 2.075$$

The ideal memory CPU with no miss is  $2.075/1.1 = 1.88$  times faster

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

# Cache Performance Example

- Suppose for the previous example we double the clock rate to 400 MHz, how much faster is this machine, assuming similar miss rate, instruction mix?
- Since memory speed is not changed, the miss penalty takes more CPU cycles:

$$\text{Miss penalty} = M = 50 \times 2 = 100 \text{ cycles.}$$

$$\text{CPI} = 1.1 + 1.3 \times .015 \times 100 = 1.1 + 1.95 = 3.05$$

$$\begin{aligned} \text{Speedup} &= (\text{CPI}_{\text{old}} \times C_{\text{old}}) / (\text{CPI}_{\text{new}} \times C_{\text{new}}) \\ &= 2.075 \times 2 / 3.05 = 1.36 \end{aligned}$$

The new machine is only 1.36 times faster rather than 2 times faster due to the increased effect of cache misses.

→ *CPUs with higher clock rate, have more cycles per cache miss and more memory impact on CPI.*



# Cache Performance

## Harvard Memory Architecture

For a CPU with separate or split level one (L1) caches for instructions and data (Harvard memory architecture) and no stalls for cache hits:

$$\text{CPUtime} = \text{Instruction count} \times \text{CPI} \times \text{Clock cycle time}$$

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

$$\text{Mem Stall cycles per instruction} =$$

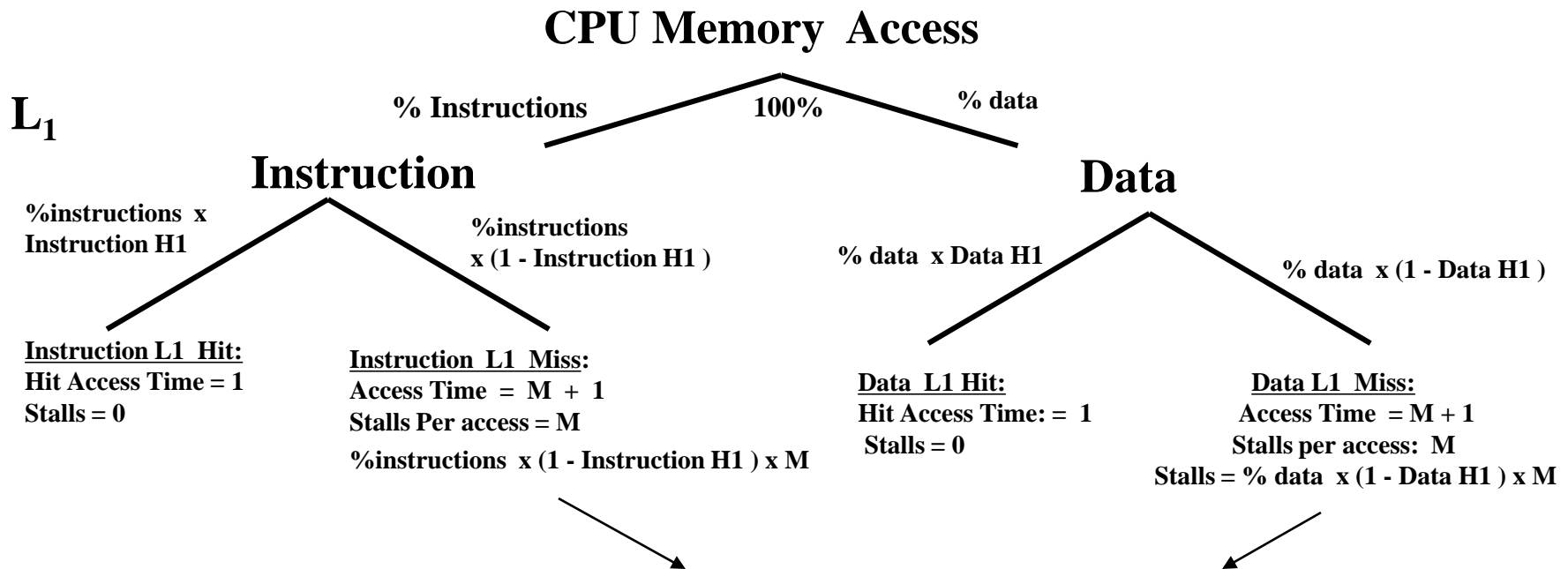
$$\text{Instruction Fetch Miss rate} \times M +$$

$$\text{Data Memory Accesses Per Instruction} \times \text{Data Miss Rate} \times M$$

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

# Memory Access Tree

## For Separate Level 1 Caches



$$\text{Stall Cycles Per Access} = \% \text{ Instructions} \times (1 - \text{Instruction H1}) \times M + \% \text{ data} \times (1 - \text{Data H1}) \times M$$

$$\text{AMAT} = 1 + \text{Stall Cycles per access}$$

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

# Cache Performance Example

- Suppose a CPU uses separate level one (L1) caches for instructions and data (Harvard memory architecture) with different miss rates for instruction and data access:
  - A cache hit incurs no stall cycles while a cache miss incurs 200 stall cycles for both memory reads and writes.
  - $CPI_{\text{execution}} = 1.1$
  - Instruction mix: 50% arith/logic, 30% load/store, 20% control
  - Assume a cache miss rate of 0.5% for instruction fetch and a cache data miss rate of 6%.
  - Find the resulting CPI using this cache? How much faster is the CPU with ideal memory?

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

$$\text{Mem Stall cycles per instruction} = \text{Instruction Fetch Miss rate} \times \text{Miss Penalty} + \\ \text{Data Memory Accesses Per Instruction} \times \text{Data Miss Rate} \times \text{Miss Penalty}$$

$$\text{Mem Stall cycles per instruction} = 0.5/100 \times 200 + 0.3 \times 6/100 \times 200 = 1 + 3.6 = 4.6$$

$$\text{Mem Stall cycles per access} = 4.6 / 1.3 = 3.5 \text{ cycles} \quad \text{AMAT} = 1 + 3.5 = 4.5 \text{ cycles}$$

$$CPI = CPI_{\text{execution}} + \text{mem stalls per instruction} = 1.1 + 4.6 = 5.7$$

The CPU with ideal cache (no misses) is  $5.7/1.1 = 5.18$  times faster

With no cache the CPI would have been  $= 1.1 + 1.3 \times 200 = 261.1 !!$

# Typical Cache Performance Data Using SPEC92

Size	Instruction cache	Data cache	Unified cache
1 KB	3.06%	24.61%	13.34%
2 KB	2.26%	20.57%	9.78%
4 KB	1.78%	15.94%	7.24%
8 KB	1.10%	10.19%	4.57%
16 KB	0.64%	6.47%	2.87%
32 KB	0.39%	4.82%	1.99%
64 KB	0.15%	3.77%	1.35%
128 KB	0.02%	2.88%	0.95%

**Miss rates for instruction, data, and unified caches of different sizes.**

# Types of Cache Misses: *The Three C's*

**1 Compulsory:** On the first access to a block; the block must be brought into the cache; also called cold start misses, or first reference misses.

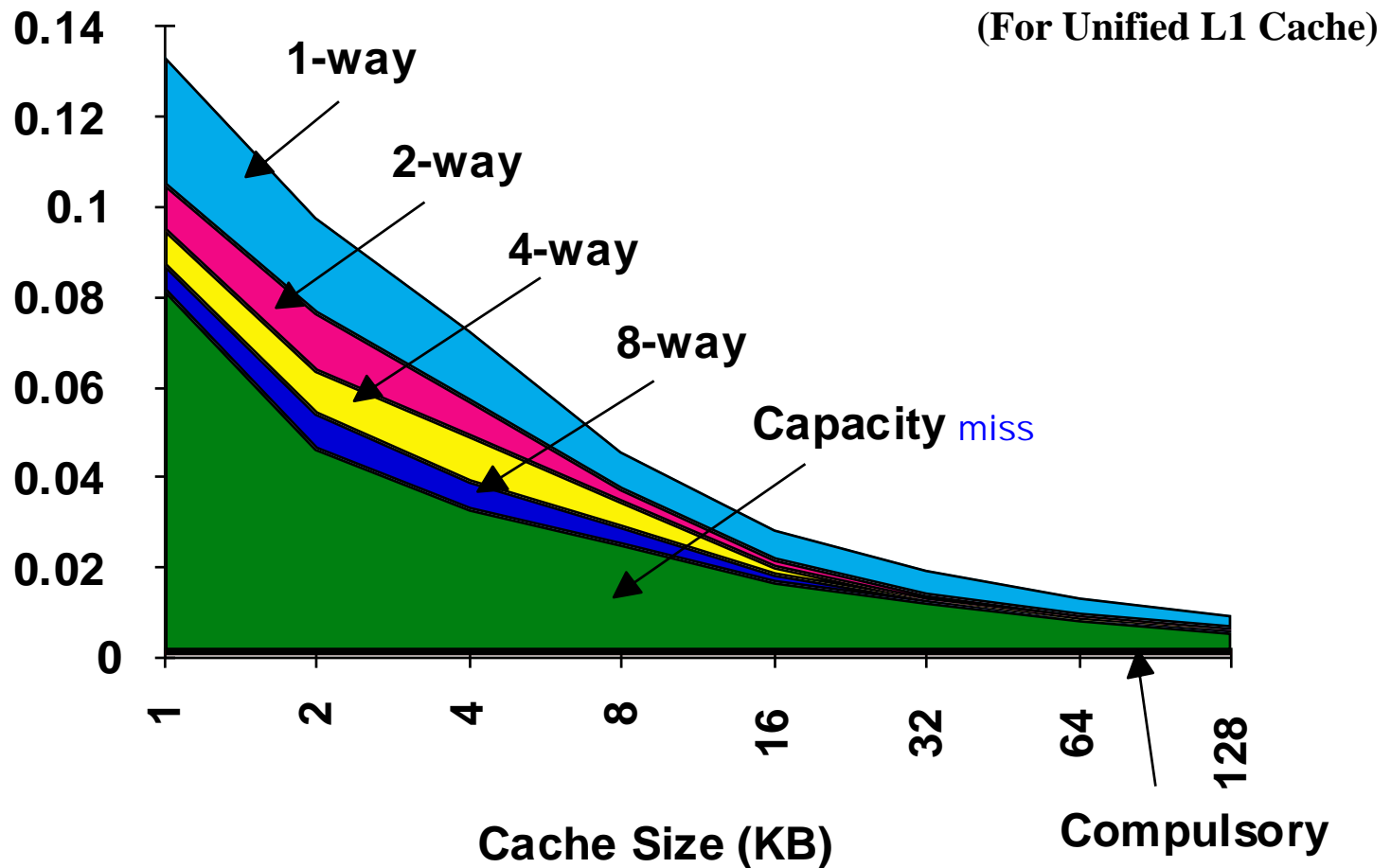
- Initially upon program startup: Miss rate ~ 100% All compulsory misses

**2 Capacity:** Occur because blocks are being discarded from cache because cache cannot contain all blocks needed for program execution (program working set is much larger than cache capacity).

**3 Conflict:** In the case of set associative or direct mapped block placement strategies, conflict misses occur when several blocks are mapped to the same set or block frame; also called collision misses or interference misses.

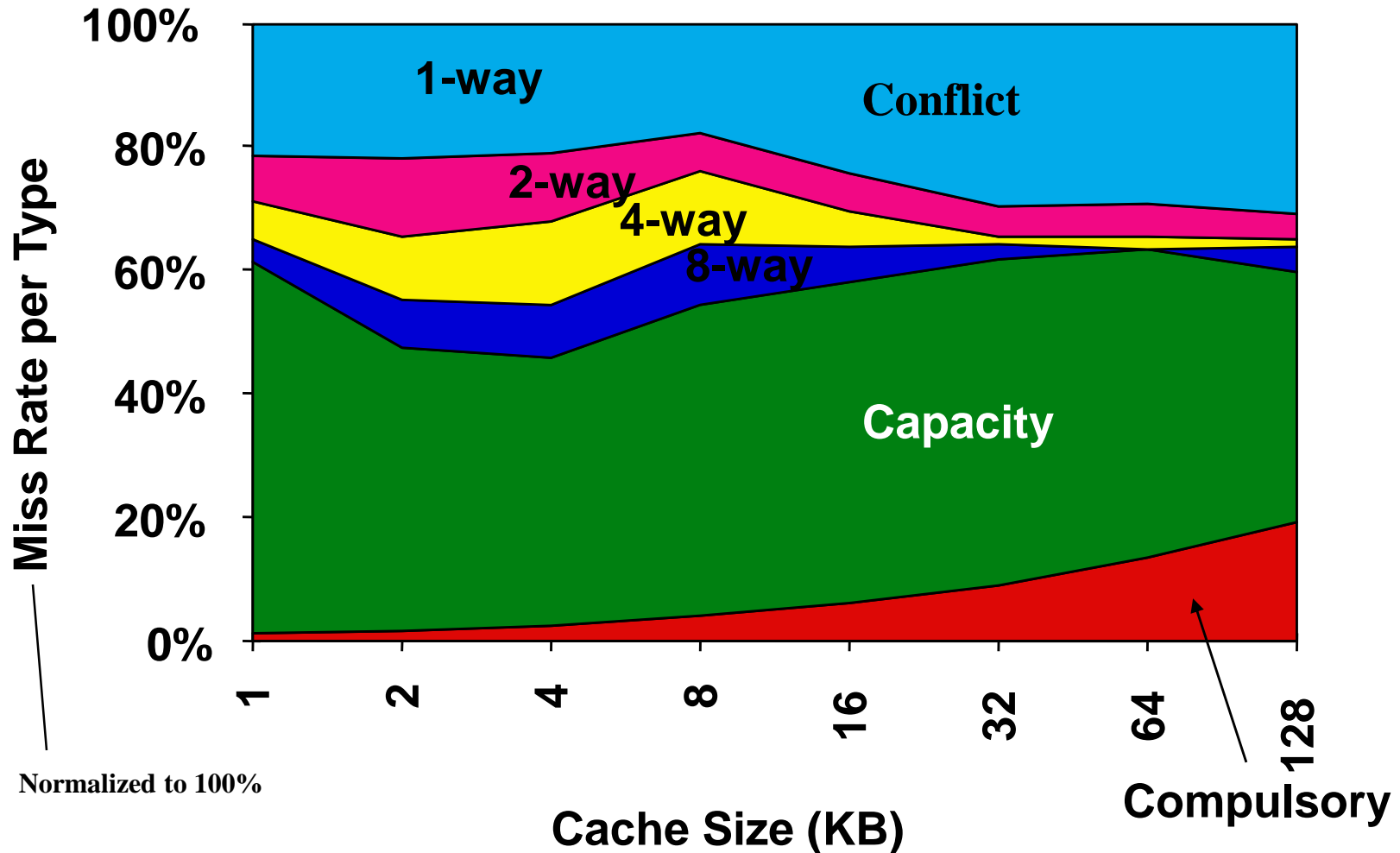
# The 3 Cs of Cache:

## Steady State Miss Rates (SPEC92)



# The 3 Cs of Cache:

## Relative Steady State Miss Rates (SPEC92)



# Cache Read/Write Operations

- Statistical data suggest that reads (*including instruction fetches*) dominate processor cache accesses (writes account for ~ 30% of data cache traffic).
- In cache reads, a cache block is read at the same time while the tag is being compared with the block address. If the read is a hit the data is passed to the CPU, otherwise ignores it.
- In cache writes, modifying the block cannot begin until the tag is checked to see if it is a hit.
- Thus for cache writes, tag checking cannot take place in parallel, and only the specific data (between 1 and 8 bytes) requested by the CPU can be modified.
- Cache can be classified according to the write and memory update strategy: write through, or write back.



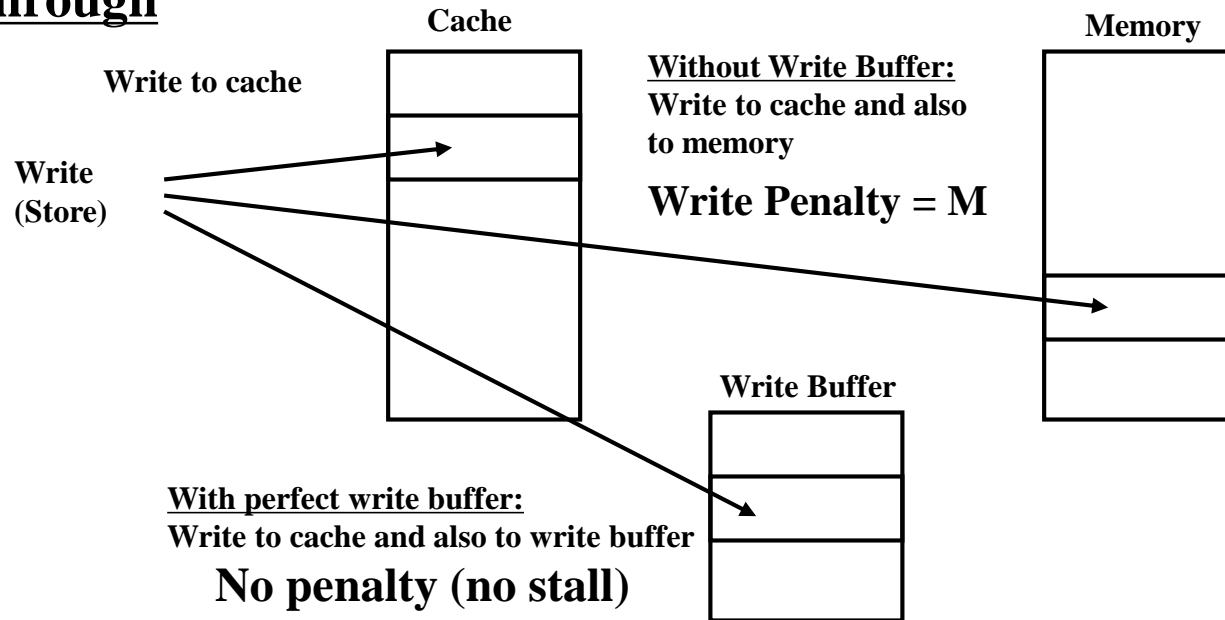
# Cache Write Strategies

- 1 **Write Through**: Data is written to both the cache block and to main memory block.
  - The lower level always has the most updated data; an important feature for I/O and multiprocessing.
  - Easier to implement than write back.
  - A write buffer is often used to reduce CPU write stall cycles while data are written to memory.
- 2 **Write Back**: Data are written or updated only to the cache block. The modified or dirty cache block is written to main memory when it's being replaced from cache.
  - Writes occur at the speed of cache
  - A status bit called a dirty or modified bit, is used to indicate whether the block was modified while in cache; if not the block is not written back to main memory when replaced.
  - **Advantage**: Uses less memory bandwidth than write through.

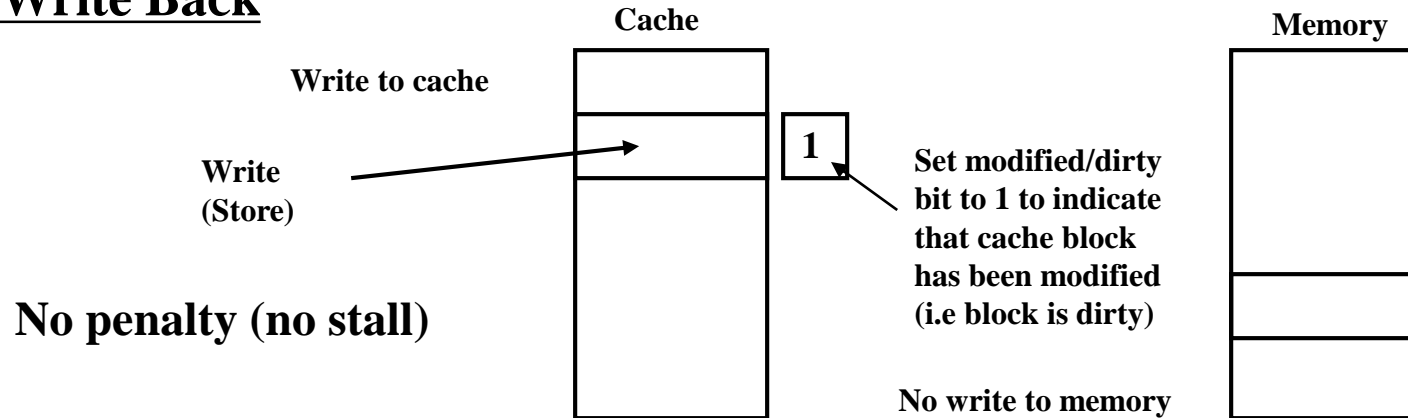
# Cache Write Strategies:

## Write Hit Operation (block to be written is in cache)

### Write Through



### Write Back



# 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:

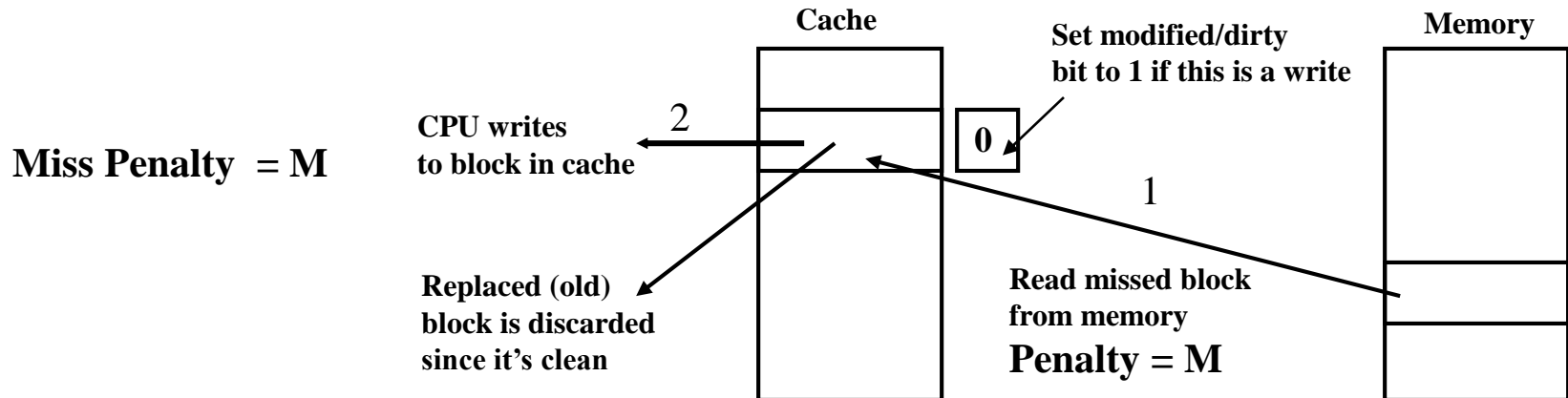
The block is modified in the lower level (lower cache level, or main memory) and not loaded into cache.

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

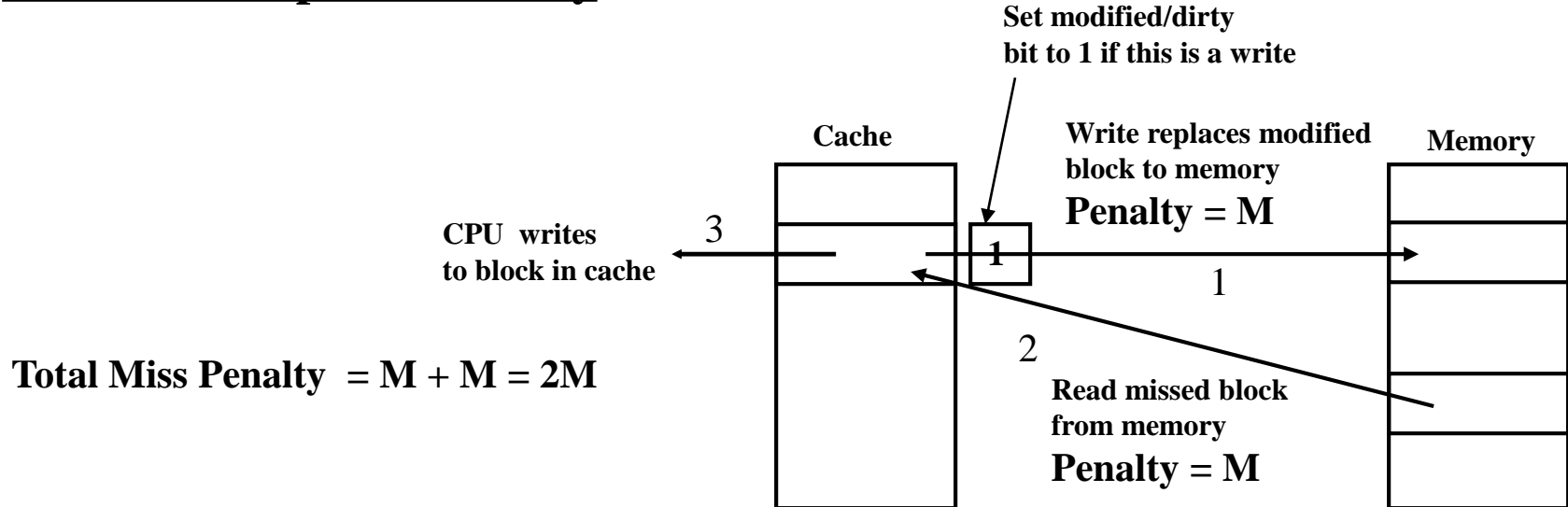
- Write back caches always use write allocate to capture subsequent writes to the block in cache.
- Write through caches usually use no-write allocate since subsequent writes still have to go to memory.

# 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**