

EECS 370: Introduction to Computer Organization Notes

Noah Peters

April 3, 2023

Abstract

Lecture notes for EECS 370 at the University of Michigan. L^AT_EXtemplate by Pingbang Hu.

Contents

3	Caches	2
3.1	Memory Hierarchy	2
3.2	Cache Basics	3
3.3	Cache Organization	5
3.4	Write Back Caches	7
3.5	Direct-Mapped Caches	9
3.6	Set Associative Caches	10
3.7	The 3 C's of Cache Misses	11
3.8	Cache Parameters vs. Miss Rate	12

Chapter 3

Caches

Lecture 17: Introduction to Caches

3.1 Memory Hierarchy

14 Mar. 12:00

We often need a lot of memory, LC2K alone can handle 2^{18} bytes of memory. We have several choices for memory:

- **SRAM:** Static Random Access Memory
 - fast: 2ns access time or faster
 - decoders are big
 - expensive, high area requirement
- **DRAM:** Dynamic Random Access Memory
 - slower: 50ns access time
 - must stall for dozens of cycles each memory load
 - less expensive
- **Flash**
 - slow: access time varies wildly
 - less expensive than DRAM
 - non-volatile
- **Disks**
 - obnoxiously slow: 3,000,000ns access time
 - dirt cheap
 - non-volatile

As seen above, there are trade-offs among each type of memory. Ideally, we would have cheap *and* fast memory. So, we can use a combination of memory types to optimize the common case via strategic **locality of reference**.

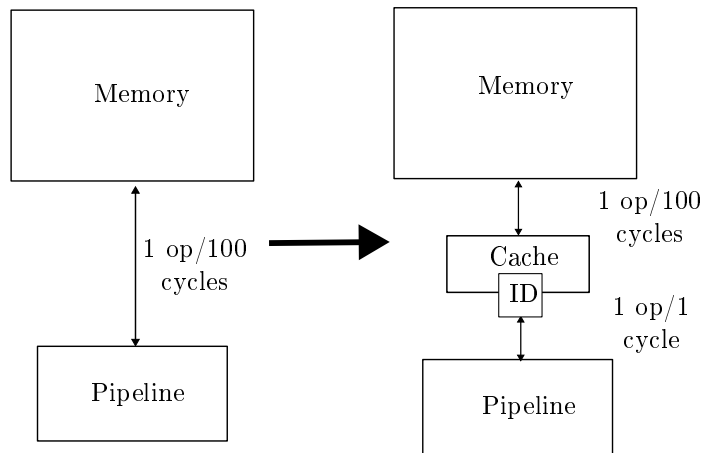


Figure 3.1: Caches Overview

Definition 3.1.1. The **architectural** view of memory is what the machine language (or programmer) sees, i.e. just a big array.

Note. Breaking up the memory system into different pieces (cache, main memory, disk) is **not architectural**. The machine language doesn't know about it.

Can use our variety of memories as follows:

- use small array of SRAM for the **cache**
- use a larger amount of DRAM for **main memory**
- use a lot of flash and/or disk for **virtual memory**

3.2 Cache Basics

Whenever memory returns data, we can store it in a cache. We'll need to store:

- the data
- a tag denoting its memory location
- a "valid" status bit

Then for our next memory access, we can first check if the tag matches the address we are attempting to access.

Definition 3.2.1 (Cache Hit). A **hit** occurs when data for a memory access is found in the cache.

Definition 3.2.2 (Cache Miss). A **miss** occurs when data for a memory access is *not* found in the cache.

Definition 3.2.3 (Hit/Miss Rate). The **hit/miss rate** is the percentage of memory accesses that hit/miss in the cache.

3.2.1 CAMs

Definition 3.2.4. CAMs: content addressable memories are akin to a set of data matching a query. Instead of an address we send a key to the memory, asking whether the key exists and, if so, what value it is associated with. Memory answers: yes/no and gives associated value (if there is one).

We apply *operations* on CAMS:

- **Search:** the primary way to access a CAM
 - send data to CAM memory
 - return "found" or "not found"
 - if found, return location of where it was found or its associated value
- **Write:** send data for CAM to remember

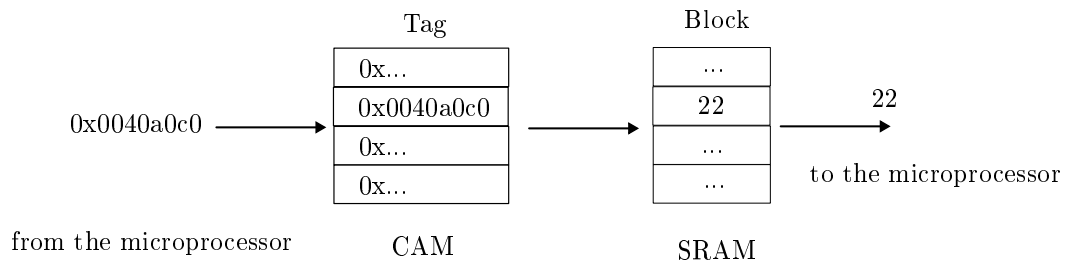
3.2.2 Cache Organization

Cache memory can copy data from any part of main memory. Cache memory has two parts:

- the **tag (CAM)**: holds the memory address
- the **block (SRAM)**: holds the memory data

A **hit** in the cache occurs when a tag match is found. The microprocessor sends the address to the CAM containing the tags and searches for the tag. If there's a search result hit, the corresponding block is forwarded to the microprocessor. If not, the address is forwarded to main memory:

Hit:



Miss:

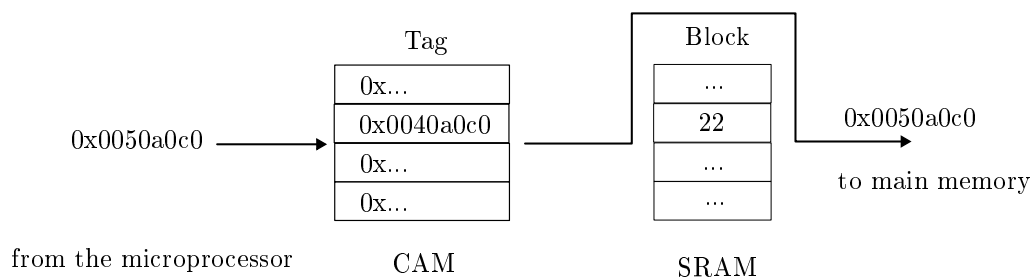


Figure 3.2: Hardware view of caches.

Problem 3.2.1. Given:

- *cache* has 1 cycle access time
- *main memory* has 100 cycle access time
- *disk* has 10,000 cycles access time

What is the average access time for 100 memory references if 90% of the cache accesses are hits

and 80% of the accesses to main memory are hits? Assume main memory access time includes tag array access to determine hit/miss.

Answer. $0.9 \cdot 1 + 0.1 \cdot (100 + 0.2 \cdot 10000) = 210.9$ ⊗

3.2.3 Cache Operation

Every cache *miss* will get the data from memory and *allocate* a cache line to put the data in (just like any CAM write). However... what do we replace in the cache? Does an optimal replacement policy exist?

Definition 3.2.5 (Temporal locality). The principle of **temporal locality** in program references says that if you access a memory location (e.g., 1000) you will be more likely to re-access that location than you will be to reference some other random location.

Remark. Locality is a property of *programs* (not hardware).

Specifically, temporal locality says that the **least recently referenced (LRU)** cache line should be *evicted* to make room for the new line. Because the re-access probability falls over time as a cache line isn't referenced, the LRU line is least likely to be re-referenced.

Definition 3.2.6 (Average Access Latency). Average Access Latency = cache latency · hit rate + memory latency ·

Lecture 18: Cache Organization: Block Size and Writes

Definition 3.2.7 (k-Way Set Associative LRU). A set associative LRU with k **ways**, has essentially k cells within each cache line's block. The LRU is set to be the *way* with the highest *count*. 20:43

16 Mar. 12:00

Each cache hit now essentially brings in two blocks of data, since each block in memory is paired up with an adjacent block. Now, we can have the LSB of our memory address correspond to either block A or block B of the retrieved block, and the remaining bits correspond to a (now-shortened) memory address.

Definition 3.2.8 (Spatial Locality). **Spatial locality** in a program says that we reference a memory location (e.g., 1000), we are more likely to reference a location near it (e.g., 1001) than some random location.

3.3 Cache Organization

We now consider the design of a cache.

3.3.1 Block Size

- choice of block size found by simulating lots of different block sizes and seeing which ones give the best performance
- most systems use a block size between 32 and 128 bytes
- *longer sizes* reduce overhead by reducing the number of CAM entries and reducing the size of each CAM entry

Problem 3.3.1. Given a cache with the following configuration:

- total size is 8 bytes

- block size is 2 bytes
 - fully associative
 - LRU replacement
 - memory address size is 16bits and is byte addressable
1. How many bits are for each tag? How many blocks in the cache?
 2. For the following reference stream, indicate whether each reference is a hit or miss: 0, 1, 3, 5, 12, 1, 2, 9, 4.
 3. What is the hit rate?
 4. How many bits are needed for storage overhead for each block?

Answer.

1. 16 bits total - 1 bit for block offset = **15 bits** for each tag. 8 bytes total / 2 bytes per block = **4 blocks**.
2.
 - 0: never in cache before → **miss**
 - 1: 0 was just cached and block size is 2 bytes → **hit**
 - 3: 2 nor 3 have been cached yet → **miss**
 - 5: 4 nor 5 have been cached yet → **miss**
 - 12: 12 nor 13 have been cached yet → **miss**
 - 1: 1 was cached and hasn't been booted from cache → **hit**
 - 2: 3 was cached recently → **hit**
 - 9: LRU was 4/5 block, so 4/5 booted from cache, 8/9 in now → **miss**
 - 4: 4/5 block just booted, so not in cache; boots 12/13 → **miss**
3. $3/9 = 0.33$
4.
 - Tag: 15 bits
 - Valid: 1 bit
 - LRU: $\log_2(4) = 2$ bits
 ⇒ total: **18 bits**

⊗

3.3.2 Stores and the Cache

Where should we write the result of a store?

- If that memory location is in the cache?
 - Send it to the cache
 - should we also send it to memory (**write-through policy**)
- If that memory location is *not* in the cache?
 - Allocate the line, i.e. put it in the cache (**allocate-on-write policy**)
 - Write it directly to memory without allocation (**no allocate-on-write policy**)

Add mem-
ory dia-
gram from
lecture

Definition 3.3.1 (Write-Through Policy). When storing a value from the processor to memory and the memory location is in the cache, we store the value in the cache and simultaneously change the value in memory.

Definition 3.3.2 (Allocate-On-Write Policy). When storing a value from the processor to memory and the memory location is *not* in the cache, we must first *read* the desired block from memory, before storing into the block in the cache and memory.

Definition 3.3.3 (No Allocate-On-Write Policy). When storing a value from the processor to memory and the memory location is *not* in the cache, we can ignore the cache and simply update memory with the given value.

Remark. This method may sometimes be less advantageous if we are going to read the value from memory soon after writing to it. However, write data streams and read data streams are often independent, so this isn't usually the case.

Lecture 19: Direct Mapped and Set Associative Caches

3.4 Write Back Caches

21 Mar. 12:00

We can design the cache to *not* write all stores to memory immediately. We can do this by keeping the most recent copy in the cache and update the memory *only when* that data is evicted from the cache.

Definition 3.4.1 (Write-Back Policy). Under this policy, when storing a value from the processor, we update the *cache* but *do not* update the value in *memory*.

We don't need to write-back all evicted lines, only those blocks that have been stored into. We can keep a **dirty bit** that *resets when the line is allocated and set when the block is stored into*. If a block is "dirty" when evicted, write its data back into memory.

V	D	LRU	Tag	Data block (bytes)
---	---	-----	-----	--------------------

Figure 3.3: Contents of a cache line.

Writes Summary

Store With No Allocate	Write-Back	Write-Through
Hit?	Write Cache	Write to Cache and <i>Mem</i>
Miss?	Write to Mem	Write to Mem
Replace Block?	If evicted block is dirty, write to Mem	Do nothing

Store With Allocate	Write-Back	Write-Through
Hit?	Write Cache	Write to Cache and <i>Mem</i>
Miss?	Read from Mem to cache, allocate to LRU block, write to cache	Read from Mem to cache, allocate to LRU block, write to cache <i>and Mem</i>
Replace Block?	If evicted block is dirty, write to Mem	Do nothing

Problem 3.4.1. Consider the following cache:

- 32-bit memory addresses
- byte addressable
- 64 KB cache
- 64 B cache block size
- **write-allocate**
- **write-back**
- *fully* associative

This cache will need 512 kilobits for the data area (64 kilobytes times 8 bits per byte). Note that here, 1 kilobyte = 1024 bytes. Besides the actually cached data, this cache will need other storage. Consider **tags**, **valid bits**, **dirty bits**, **bits to track LRU**, etc.

How many additional bits (not counting data) will be needed to implement this cache?

Answer. We have 2^{10} blocks:

$$64 \text{ KB} = 2^{16} \text{ b} \rightarrow \frac{2^{16} \text{ b}}{2^6 \text{ b / block}} = 2^{10} \text{ blocks.}$$

The **tag's** size:

$$32 \text{ b for mem addresses} - 6 \text{ b for block offset} = 26 \text{ b for tag.}$$

The **valid** and **dirty bits** only need 1 bit each. To keep track of LRU we can calculate the number of bits needed as follows:

$$2^{16}/2^6 = 2^{10} \text{ blocks in the cache} \rightarrow 10 \text{ bits needed to keep track of LRU}$$

So, the total number of additional bits needed for the cache is

$$26 + 1 + 1 + 10 = \mathbf{38 \text{ bits}} \text{ (per block).}$$

⊛

Problem 3.4.2. Suppose that accessing a cache takes **10 ns** while accessing main memory, and in the case of cache-miss it takes **100 ns**.

- a What is the average memory access time if the cache hit rate is 97%?
- b If the cache size is increased, causing the hit rate to rise by 1% and the time for accessing the cache by 2 ns. Will this improve performance?

Answer.

a

$$10 \text{ ns} + 0.03 \cdot 100 \text{ ns} = 13 \text{ ns.}$$

b

$$12 \text{ ns} + 0.02 \cdot 100 \text{ ns} = 14 \text{ ns.}$$

⊛

3.5 Direct-Mapped Caches

A block can go to *any* location. This means that when we are checking tags, we *check every cache tag* to determine whether the data is in the cache. This parallel approach is what is used for **fully associative caches**, which we have studied so far. However, this method is *slow*.

We can redesign the cache to eliminate the requirement for parallel tag lookups.

Definition 3.5.1 (Direct-Mapped Caches). Direct-mapped caches partition memory into as many regions as there are cache lines. Each memory region maps to a **single cache line** in which data can be placed. You then only need to **check a single tag**: the one associated with the *region the reference is located in*.

Remark. The memory regions that map to the same cache line are placed as far apart in memory as possible, so that blocks that are spatially close are not competing for the same cache line.

Since two blocks in memory that map to the same index in the cache cannot be present in the cache at the same time, a 0% hit rate is possible if more than one block accessed in an interleaved manner map to the same index.

Problem 3.5.1. How many tag bits are required for the following cache specs? What are the overheads of this cache?

- 32-bit address, byte addressed
- 128 byte block size
- direct-mapped 32k cache
- write-back

Answer. The number of bits required for the line index is

$$15 - 7 = 8. (2^{15} = 32 \text{ kb cache})$$

The tag size is

$$32 \text{ b total} - 7 \text{ b for block offset} - 8 \text{ b for line index} = \mathbf{17 \text{ b for tag.}}$$

So the overhead in this cache is

$$\begin{aligned} 17 + 1 \text{ b for valid bit} + 1 \text{ b for dirty bit} &= 19 \text{ b / cache line} \\ 19 \text{ b / line} \cdot 256 \text{ lines} &= 4864 \text{ bits} \\ 4864 \text{ b} / 32 \text{ KB} &= \mathbf{1.9\% \text{ overhead.}} \end{aligned}$$

⊗

Lecture 20: Set-Associative Caches and 3 C's

Unlike fully-associative caches, in direct-mapped caches the tag array (CAM) doesn't have to be searched *before* the data array (SRAM). Instead, we can do a **direct lookup** and search the tag array and data array in **parallel**, both of which are faster for cache lookups. The *downside* to direct-mapped caches, however, is that there is an increased chance for collisions in the cache.

23 Mar. 12:00

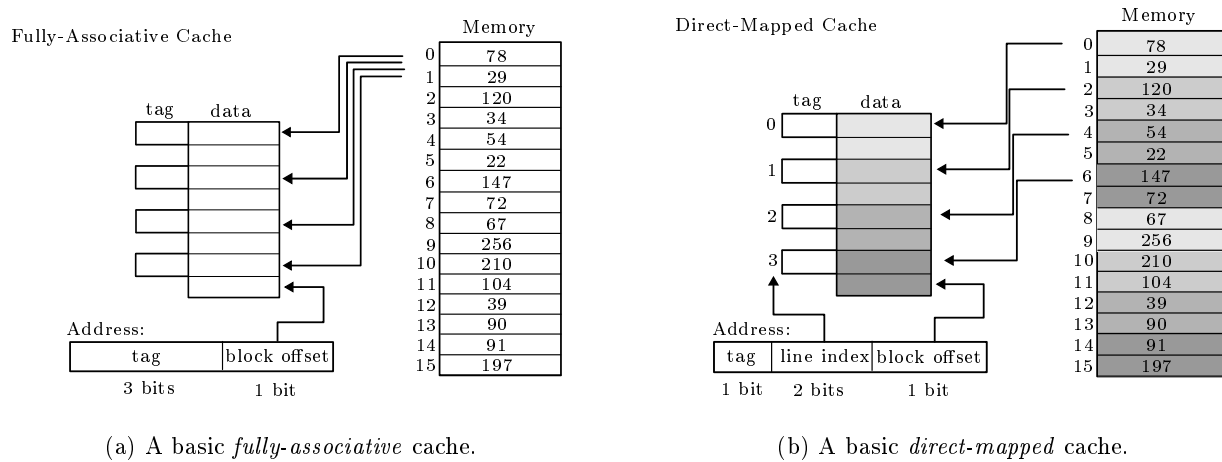


Figure 3.4: Fully-associative and direct-mapped caches.

3.6 Set Associative Caches

We can achieve the advantages each the fully-associative and direct-mapped cache structures using **set associative caches**.

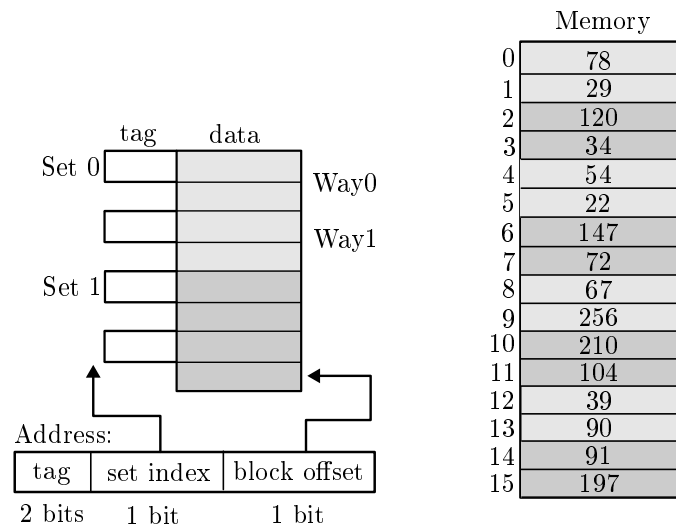


Figure 3.5: A basic set-associative cache.

Definition 3.6.1 (Set Associative Cache). Set associative caches:

- Partition memory into regions (like direct-mapped but fewer partitions)
- Associate a region to a set of cache lines; check tags for all lines in a set to determine a hit
- Treat each line in a set like a small fully associative cache (with LRU policy)

Problem 3.6.1. For a 32-bit address and 16 KB cache with 64-byte blocks, show the breakdown of the address for the following cache configurations:

- Fully Associative Cache
- 4-Way Set Associative Cache
- Direct-Mapped Cache

Answer. Since the cache size is $16 \text{ KB} = 2^4 \cdot 2^{10} = 2^{14}$ bytes and there are 64 bytes per block we have

$$2^{14}/2^6 = 2^8 \text{ blocks.}$$

a For the **fully associative** cache:

64 byte blocks \rightarrow 6 bits for block offset

$$32 \text{ b} - 6 \text{ b} = 26 \text{ bits for tag}$$

b For the **4-way associative** cache:

64 byte blocks \rightarrow 6 bits for block offset

$$2^8 \text{ blocks} / 4 \text{ blocks per set} = 2^6 \text{ sets} \rightarrow 6 \text{ bits for set index}$$

$$32 \text{ b} - 12 \text{ b} = 20 \text{ bits for tag}$$

c For the **direct-mapped** cache:

64 byte blocks \rightarrow 6 bits for block offset

2^8 different blocks \rightarrow 8 bits for line index

$$32 \text{ b} - 6 - 8 = 18 \text{ bits for tag}$$

⊛

Remark. Note that we've used *block sizes*, *number of sets*, *number of ways* that are all be power of 2. This allows us to properly space blocks, sets, and data, etc.

3.7 The 3 C's of Cache Misses

Definition 3.7.1 (Compulsory Miss). A **compulsory** miss, also called a "cold start" miss, occurs the first time there is a reference to any block, which always results in a miss.

Definition 3.7.2 (Capacity Miss). A **capacity** miss occurs when the cache is too small to hold all the data, and a hit would have occurred if the cache were large enough.

Definition 3.7.3 (Conflict Miss). A **conflict** miss occurs when the cache is not associated enough. One set is getting a disproportionate amount of accesses, where a hit would have occurred with a fully associative cache.

3.7.1 Classifying Cache Misses

To classify cache misses we can use different forms of simulation:

- Simulate with a cache of unlimited size (cache size = memory size) \rightarrow any misses must be *compulsory* misses
- Simulate again with a fully associative cache of intended size \rightarrow any new misses must be *capacity* misses
- Simulate a third time, with the actually intended cache \rightarrow any *new* misses must be *conflict* misses

3.7.2 Fixing Cache Misses

We have various methods of reducing each type of cache miss:

- *Compulsory*: increase block size (and thus reduce total number of blocks)
- *Capacity*: build bigger cache
- *Conflict*: increase associativity

Problem 3.7.1. Consider a cache with the following configuration: write-allocate, total size of 64 bytes, block size is 16 bytes and 2-way associative, memory address size is 16 bits and byte-addressable, replacement policy is LRU, and cache is empty at the start.

For the following memory accesses, indicate whether the reference is a hit or miss, and the type of miss:

0x00, 0x14, 0x27, 0x08, 0x38, 0x4A, 0x18, 0x27, 0x0F, 0x40

Answer. Note that the block size is 16 bytes, so FA and SA have 4 blocks. SA is 2-way.

Address	Infinite	FA	SA	3 Cs
0x00	M	M	M	Compulsory
0x14	M	M	M	Compulsory
0x27	M	M	M	Compulsory
0x08	H	H	H	N/A
0x38	M	M	M	Compulsory
0x4A	M	M	M	Compulsory
0x18	H	M	H	N/A
0x27	H	M	M	Capacity
0x0F	H	M	M	Capacity
0x40	H	H	M	Conflict

⊗

3.8 Cache Parameters vs. Miss Rate

3.8.1 Cache Size

As cache total data (not including tag) capacity increases, temporal locality can be used better. However, it's not *always* better. Too large of a cache adversely affects hit and miss latency, too small of a cache doesn't exploit temporal locality well. The **working set** (i.e. the whole set of data executing application references within a time interval) size is constant, so expanding cache size beyond that isn't necessarily advantageous.

3.8.2 Block Size

Block size is the data that is associated with an address tag. Too small blocks don't exploit spatial locality well and have larger tag overhead. Too large blocks have too few total number of blocks, so we're less likely to transfer useless data.

3.8.3 Associativity

How many blocks can map to the same index? Larger associativity causes lower miss rates, less variation among programs, but diminishing returns. Smaller associativity causes lower costs, faster hit time.

Lecture 21: Cache Wrap-Up

Problem 3.8.1. The *grinder* app running on LC2K with full data forwarding and all branches predicted not-taken has the following frequencies:

45% R-type, 20% branches 15% loads, 20% stores

In *grinder*, 40% of branches are taken and 50% of LWs are followed by an immediate use. The I-cache has a miss rate of 3% and the D-cache has a miss rate of 6% (no overlapping of misses). On a miss, the main memory is accessed and has a latency of 100 ns. The clock frequency is 500 MHz.

What is the CPI of *grinder* on the LC2K?

Answer. Since the clock frequency is 500 MHz, we have a 2 ns cycle time. So the stalls per cache miss have a length of

$$100 \text{ ns} / 2 \text{ ns} = 50 \text{ cycles.}$$

So the CPI will equal:

$$\begin{aligned} CPI &= 1 + \text{data hazard stalls} + \text{ctrl hazard stalls} + \text{I-Cache stalls} + \text{D-Cache stalls} \\ &= 1 + 0.15 \cdot 0.5 \cdot 1 + 0.2 \cdot 0.4 \cdot 3 + 1 \cdot 0.03 \cdot 50 + 0.35 \cdot 0.06 \cdot 50 = \mathbf{3.865}. \end{aligned}$$

⊛

Problem 3.8.2. Say you have the following:

- *program*: generates 2 billion loads, 1 billion stores, each 4 bytes in size
- *cache*: a 32-byte block which gets a 95% hit rate on the program

How many bytes of memory would be read and written if:

- a We had no cache?
- b We had a write-through cache with a no-write allocate policy?
- c We had a write-back cache with a write-allocate policy? (Assume 25% of misses result in dirty eviction)

Answer.

- a *No cache*: All stores go to memory and are 4 bytes each, so 1 billion · 4 bytes = 4 billion bytes for writes. Similarly, there are 4 bytes per read, so there are 2 billion · 4 bytes = 8 billion bytes for reads.
- b *Write-through, no allocate*: No allocate means we bypass cache on a miss, and write-through means we write to cache and memory on hit. So, again all stores will go to memory and are still 4 bytes each, so there are **4 billion bytes of writes**. Only loads that miss in the cache go to memory, but they read the full cache block. So there are 2 billion · 0.05 · 32 bytes = **3.2 billion bytes of reads**.
- c *Write-back, write-allocate*: In this case, store misses result in a cache block being read, causing 1 billion stores · 0.05 · 32 bytes = 1.6 billion bytes to be read due to store misses. Similarly, *load* misses result in cache block being read, causing an additional 2 billion loads · 0.05 · 32 bytes = 3.2 billion bytes to be read. Thus, there are **4.8 billion bytes read**. We only write to memory when there are dirty evictions, which can be done by both loads and stores. Since there are 0.05 · 1 billion stores · 32 bytes · 25% = 0.4 billion bytes of writes due to store misses. There are 2 billion loads · 0.05 · 32 bytes · 25% = 0.8 billion bytes of writing due to load misses, resulting in **1.2 billion bytes of writes** in total.

Problem 3.8.3. Given a 200 MHz processor with 8 KB instruction and data caches and a with memory access latency of 20 cycles. Both caches are 2-way associative. A program running on this processor has a 95% icache hit rate and a 90% dcache hit rate. On average, 30% of the instructions are loads or stores. The CPI of this system, if caches were ideal would be 1.

Suppose you have the following two options for the next generation processor, which do you pick?

1. Double the clock frequency → assume this will increase your memory latency to 40 cycles, and the base CPI of 1 can still be achieved after this change
2. Double the size of your caches → this will increase the instruction cache hit rate to 98% and the data cache hit rate to 95%; assume the hit latency is still 1 cycle

Answer. See end of lecture.

Appendix