CS 211: Intro to Computer Architecture 13.2: Memory Access Patterns and Caching

Minesh Patel Spring 2025 – Thursday 24 April

Announcements

- Ongoing
 - WA9: due Friday (April 18) @ 11:59 pm
 - PA5: due Monday (May 5) @ 11:59 pm
- Upcoming
 - WA10: TBA: planned for Friday
 - Final Exam: May 14th

Our Mental Model (So Far)



- In reality, we have **different types of memory**
- Today: registers vs. caches vs. main memory

Harris & Harris, "Digital Design and Computer Architecture: RISC-V Edition"



•Why Is Main Memory Slow?

•Caching and the Memory Hierarchy

Direct-Mapped Caches
Where to Place Data
How to Query the Cache

Semiconductor Memories



FIGURE 12.1 Categories of memory arrays

Weste & Harris, "CMOS VLSI Design", 4/E

Main Memory vs. Registers

Registers (Flip-Flops)

Built from digital logic

Main Memory (DRAM)

Built from analog circuits



Main Memory vs. Registers





Registers (Flip-Flops) Built from digital logic Main Memory (DRAM) Built from analog circuits

Cost	Expensive	
Speed	Fast	
Physical Size	Big	

Cheap Slow Small

Even If I Were a Billionaire....

• I still couldn't afford to keep all data in CPU registers



Beyond Cost: Access Latency

- Fundamentally: **bigger = slower**
 - Electrical signals travel further
 - More wire resistance
 - ...
- Low-capacity memories are fastest



Beyond Cost: Density

Static RAM (Registers, Caches)



Dynamic RAM (Main Memory)

Built from analog circuits



Prof. Mason, ECE 410, Lecture 13.1

Weste & Harris, "CMOS VLSI Design", 4/E **10**

Beyond Cost: Energy

Energy Cost of Various CPU Operations



Horowitz, Mark. "Computing's Energy Problem (and What We Can Do About It)," ISSCC, 2014.

Cost-Performance Tradeoff



• There are **other concerns** too (e.g., energy, reliability, etc...)

Harris & Harris, "Digital Design and Computer Architecture: RISC-V Edition"

It's Just Hard to Improve Memory



Hennessy & Patterson, "A Quantitative Approach to Computer Architecture," 6/E., Figure 2.2



•Why Is Main Memory Slow?

Caching and the Memory Hierarchy

Direct-Mapped Caches Where to Place Data How to Query the Cache

Main Memory Kills Performance

• The CPU wastes lots of time just moving data

Instruction type



Getting the Best of Both Worlds

• Goal: Give the illusion that main memory is large but fast



- If we cache the right data, the CPU never goes to main memory
 - Common case: accesses are at the speed of cache
 - Worst case: data not in the cache

The Right Data to Cache



• A: Data that we know we will need soon (prediction problem!)

- Cache hit: data is in the cache fast
- Cache miss: data is NOT in the cache slow

Choosing What to Cache

• The answer comes from studying application behavior



func:

.L2:

slli a1, a1, 3 add a1, a0, a1

Memory Access Pattern

read(A[0])
write(A[0])
read(A[1])
write(A[1])
Observation 2: Adjacency
write(A[1])

read(A[N-1])
write(A[N-1])

Benchmarking ilab2

for(uint64_t i = 0; i < array_size; i += stride)
 temp = arr[i];</pre>



li	a4, 0
.loop:	
bge	a4, a2, .done
ld	t1, 0(t0)
add	a4, a4, a3
add	t0, t0, a3
j	.loop
. done:	
•••	

Example: Matrix Multiply (C = A*B)



(a) Version *i j k*

CS:APP 3/E, "Section 6.6.2"

The Memory Hierarchy

- Key idea: Have successively larger caches
 - "Hot data" close to CPU (closer the better)
 - "Cold data" far from CPU



Hennessy and Patterson, "Computer Architecture" 6/E.

Hit/Miss Performance

• Performance depends on whether there is a cache hit or miss



- •L1 hit:
- •L1 miss + L2 hit:
- •L1 miss + L2 miss + L3 hit:
- •L1 miss + L2 miss + L3 miss + Memory:

Testing Memory/Cache Latency



https://www.techpowerup.com/forums/threads/share-your-aida-64-cache-and-memory-benchmark-here.186338/page-99

Cache Design Questions

- Which data should we cache?
- How do we access the cache?
- Where in the cache should that data go?
- •What happens when the cache is full?
- •What if we **modify** data in the cache?



•Why Is Main Memory Slow?

Caching and the Memory Hierarchy

Direct-Mapped Caches
Where to Place Data
How to Query the Cache

Cache Blocks

- Caches store **"blocks" of data** (usually 64 bytes)
- Think about it like a smaller version of main memory



Main Memory 2⁶⁴ / 64 blocks



Types of Caches

• Q: Suppose we issue an 1d (8B): which block does it go in?



Fully Associative: Anywhere

Set-Associative: Specific places

Direct Mapped: Exactly one place

Today

Cache Organizations: Terminology

• Rows (sets) and columns (ways) of blocks



Set-Associative



Fully-Associative



Toward a Direct-Mapped Cache



Set 2

Set 3



File

- •Where to place data?
- •How to query the cache?

Memory



•Why Is Main Memory Slow?

•Caching and the Memory Hierarchy

Direct-Mapped Caches
Where to Place Data
How to Query the Cache

Mapping Data to a Direct-Mapped Cache

• Need to map every memory location to a specific set



Set Index Hash Function

We could use something fancy (lots of research on this!)
Simplest: just use bits of the address



Visualizing the Cache Mapping



Resolving Hash Collisions

Solution: cache both data and address

	"Tag"	Set Index	Byte-in-Block
64-bit address	a ₆₃ a ₁₀ a ₉	a ₈ a ₇ a ₆	$a_5a_4a_3a_2a_1a_0$

Direct-Mapped Cache



Must compare tag bits

on every cache access to check for collision



•Why Is Main Memory Slow?

•Caching and the Memory Hierarchy

Direct-Mapped Caches
Where to Place Data
How to Query the Cache

Querying the Cache

- The cache must **track the current state** of every block
 - INVALID: no data currently cached (tag is meaningless)
 - CLEAN: data is cached and unmodified
 - **DIRTY:** data is cached and modified by the CPU (doesn't match memory)

Data	Tag	State
	Data	Data Tag

Direct-Mapped Cache

Querying the Cache: Cache Controller



Querying the Cache: STORE Instruction

Direct-Mapped Cache



func:	
//	a0 = 0x100
//	a1 = 0xff
sd	a1, 0(a0)

Example: Accessing an Array

Direct-Mapped Cache



```
void sum(uint64_t a[10])
{
    // a = 0x100
    uint64_t s = 0;
    for(int i = 0; i < 10; i++)
        s += a[i];
    return s;
}</pre>
```

Example: Accessing an Array

 Register
 Set 0
 Image: Set 2
 Image: Set 3
 Image: Set

<pre>void sum(uint64_t a[10]) {</pre>
l
// a = 0x100
uint64 + s = 0
$\operatorname{diff}(\mathbf{O} + \mathbf{C} \mathbf{O} + \mathbf{O})$
<pre>for(int i = 0; i < 10; i++)</pre>
c = c
S ∓- a[⊥],
return s:
}
J

a[0]	=	1_ <mark>00</mark> 00_0000
a[1]	=	1 <u>00</u> 00_1000
a[2]	=	1 <u>00</u> 01_0000
a[3]	=	1 <u>00</u> 01_1000
a[4]	=	1 <u>00</u> 10_0000
a[<mark>5</mark>]	=	1 <u>00</u> 10_1000
a[<mark>6</mark>]	=	1 <u>00</u> 11_0000
a[7]	=	1 <u>00</u> 11_1000
a[<mark>8</mark>]	=	1_ <mark>01</mark> 00_0000
a[9]	=	1_0100_1000

<pre>read(a[0])</pre>	=	MISS
<pre>read(a[1])</pre>	=	HIT
<pre>read(a[2])</pre>	=	HIT
<pre>read(a[3])</pre>	=	HIT
<pre>read(a[4])</pre>	=	HIT
<pre>read(a[5])</pre>	=	HIT
<pre>read(a[6])</pre>	=	HIT
<pre>read(a[7])</pre>	=	HIT
<pre>read(a[8])</pre>	=	MISS
<pre>read(a[9])</pre>	=	HIT

CS 211: Intro to Computer Architecture 13.2: Memory Access Patterns and Caching

Minesh Patel Spring 2025 – Thursday 24 April