

Overview

- Notes

- ▣ Homework 9 is due tonight

- Verify your submitted file before midnight

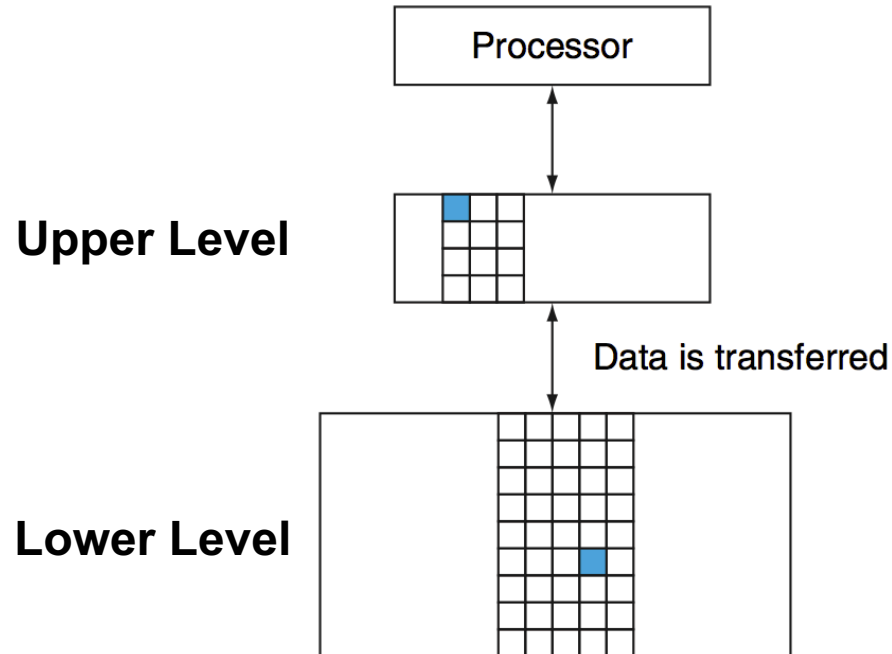
- This lecture

- ▣ Cache

Recall: Memory Hierarchy

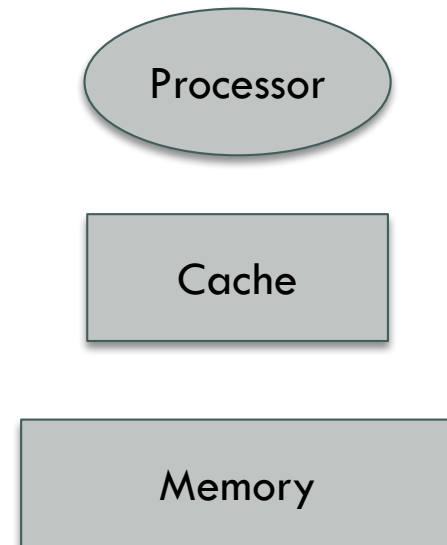
- The basic structure of a memory hierarchy.
- Multiple levels of the memory

Idea: keep important data closer to processor.



Cache Architecture

- Design principles
 - ▣ Temporal locality: if you used some data recently, you will likely use it again
 - ▣ Spatial locality: if you used some data recently, you will likely access its neighbors
- Cache terminology
 - ▣ Access time
 - ▣ Hit vs. miss
 - ▣ Miss penalty



Cache Terminology

- Block (cache line): unit of data access
- Hit: accessed data found at current level
 - ▣ *hit rate: fraction of accesses that finds the data*
 - ▣ *hit time: time to access data on a hit*
- Miss: accessed data NOT found at current level
 - ▣ miss rate: $1 - \text{hit rate}$
 - ▣ miss penalty: time to get block from lower level

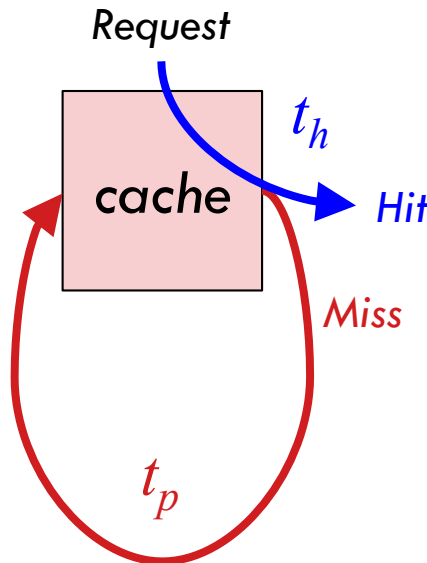
hit time \ll miss penalty

Cache Performance

□ Average Memory Access Time (AMAT)

Outcome	Rate	Access Time
Hit	r_h	t_h
Miss	r_m	$t_h + t_p$

$$AMAT = r_h t_h + r_m (t_h + t_p)$$
$$r_h = 1 - r_m$$



$$AMAT = t_h + r_m t_p$$

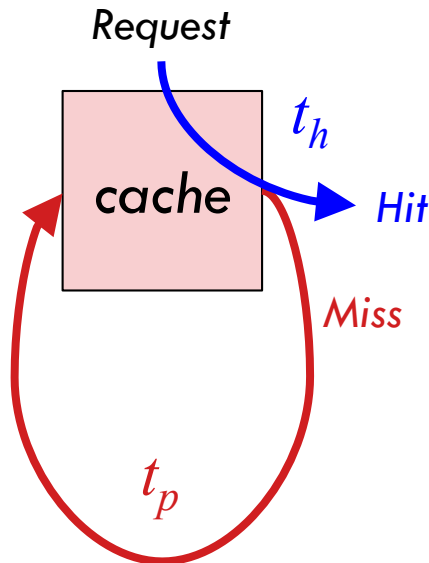
problem: hit rate is 90%; hit time is 2 cycles;
and accessing the lower level takes 200 cycles;
find the average memory access time?

Cache Performance

□ Average Memory Access Time (AMAT)

Outcome	Rate	Access Time
Hit	r_h	t_h
Miss	r_m	$t_h + t_p$

$$AMAT = r_h t_h + r_m (t_h + t_p)$$
$$r_h = 1 - r_m$$



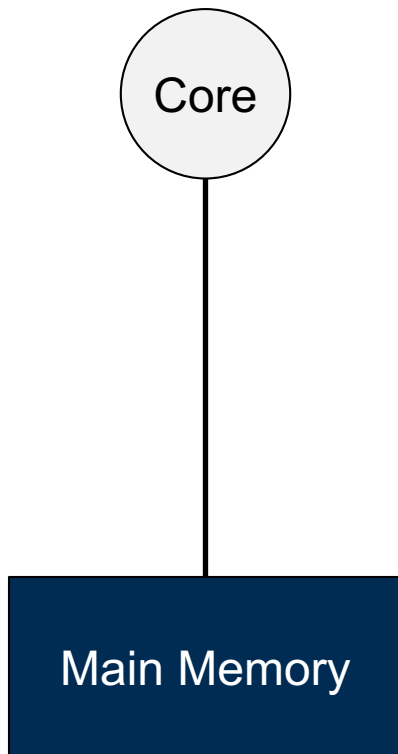
$$AMAT = t_h + r_m t_p$$

problem: hit rate is 90%; hit time is 2 cycles;
and accessing the lower level takes 200 cycles;
find the average memory access time?

$$AMAT = 2 + 0.1 \times 200 = 22 \text{ cycles}$$

Summary: Cache Performance

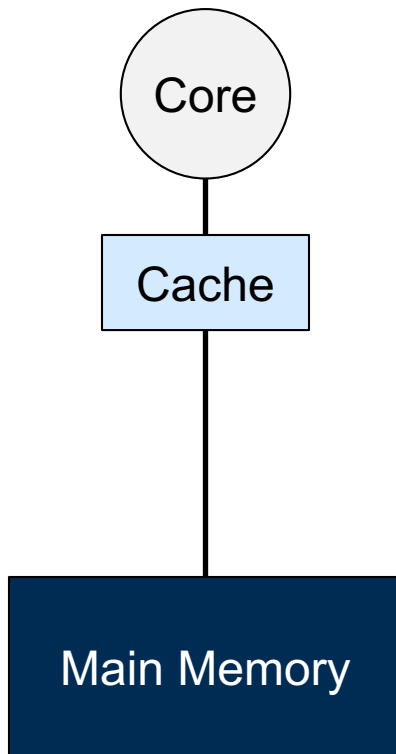
- Bridging the processor-memory performance gap



Main memory access time: 300 cycles

Summary: Cache Performance

- Bridging the processor-memory performance gap



Main memory access time: 300 cycles

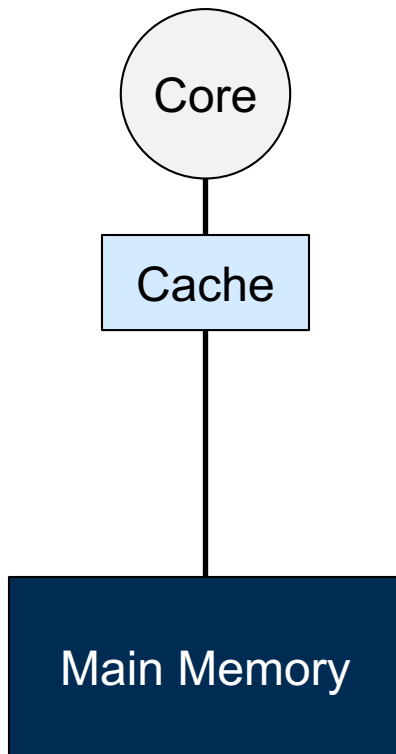
Cache

- L1: 2 cycles hit time; 60% hit rate

What is the average mem access time?

Summary: Cache Performance

- Bridging the processor-memory performance gap



Main memory access time: 300 cycles

Cache

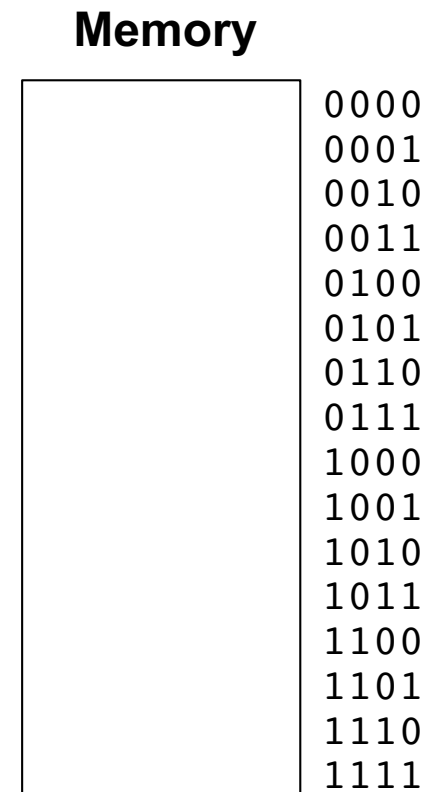
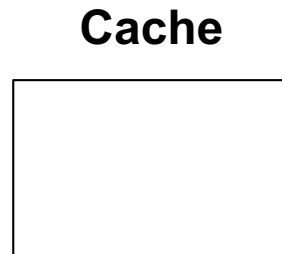
- L1: 2 cycles hit time; 60% hit rate

What is the average mem access time?

$$\begin{aligned} AMAT &= t_h + r_m t_p \\ &= 2 + 0.4 \times 300 \\ &= 122 \end{aligned}$$

Cache Addressing

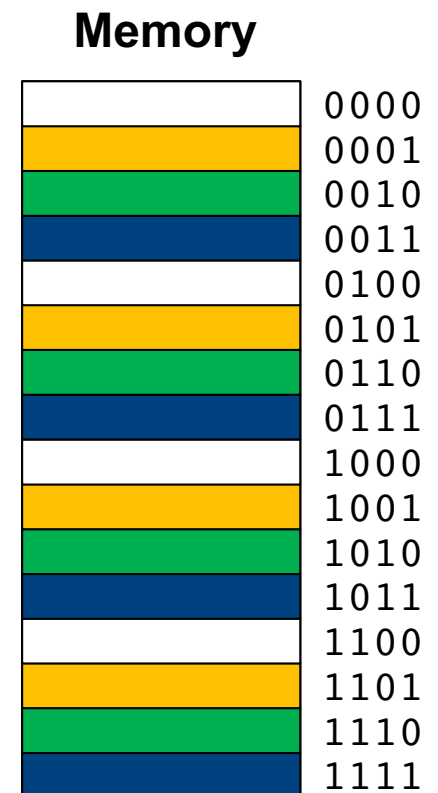
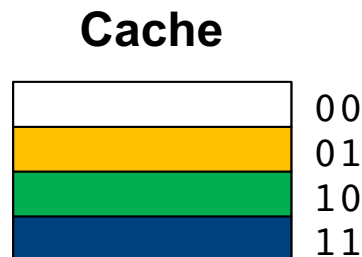
- Instead of specifying cache address we specify main memory address
- Simplest: direct-mapped cache



Cache Addressing

- Instead of specifying cache address we specify main memory address
- Simplest: direct-mapped cache

Note: each memory address maps to a single cache location determined by modulo hashing

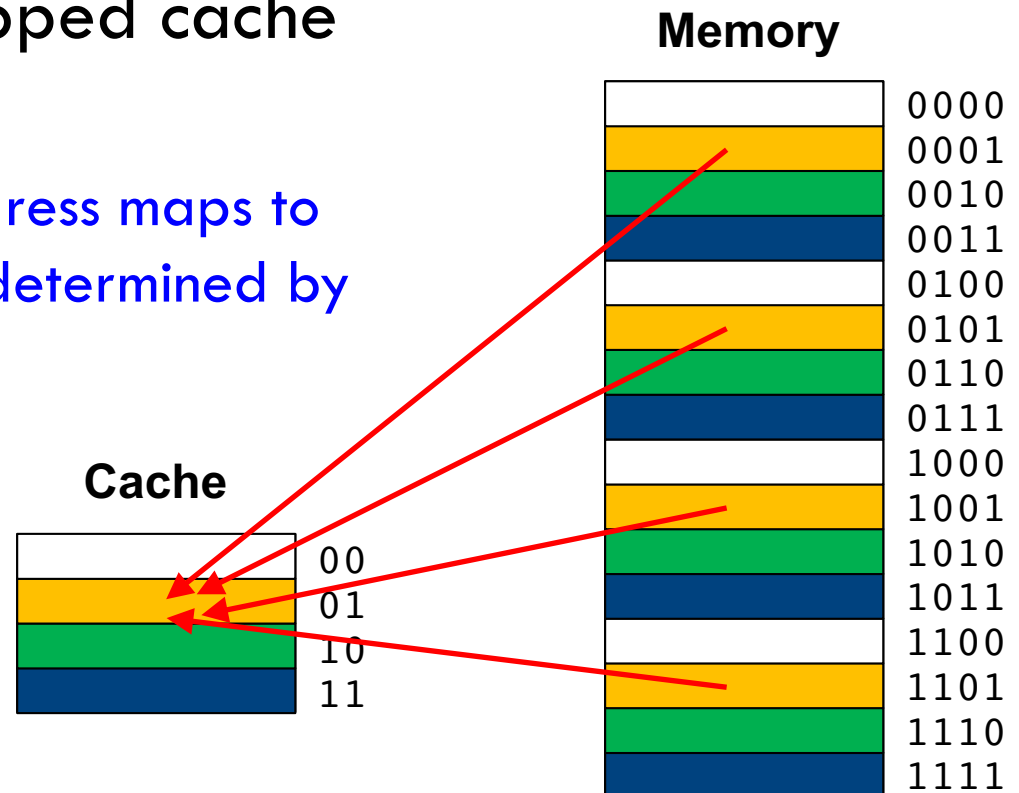


Cache Addressing

- Instead of specifying cache address we specify main memory address
- Simplest: direct-mapped cache

Note: each memory address maps to a single cache location determined by modulo hashing

How to exactly specify which blocks are in the cache?



Direct-Mapped Lookup

- ❑ Byte offset: to select the requested byte
- ❑ Tag: to maintain the address
- ❑ Valid flag (v): whether content is meaningful
- ❑ Data and tag are always accessed

