

MK Computer Architecture
A Quantitative Approach, Fifth Edition

Chapter 2
Memory Hierarchy Design

MK Copyright © 2012, Elsevier Inc. All rights reserved. 1

Introduction

- Programmers want unlimited amounts of memory with low latency
- Fast memory technology is more expensive per bit than slower memory
- Solution: organize memory system into a hierarchy
 - Entire addressable memory space available in largest, slowest memory
 - Incrementally smaller and faster memories, each containing a subset of the memory below it, proceed in steps up toward the processor
- Temporal and spatial locality insures that nearly all references can be found in smaller memories
 - Gives the illusion of a large, fast memory being presented to the processor

MK Copyright © 2012, Elsevier Inc. All rights reserved. 2

Memory Hierarchy

Component	Size	Speed
Register reference	1000 bytes	300 ps
Level 1 Cache reference	64 KB	1 ns
Level 2 Cache reference	256 KB	3-10 ns
Level 3 Cache reference	2-4 MB	10-20 ns
Memory reference	4-16 GB	50-100 ns
Disk memory reference	4-16 TB	5-10 ms

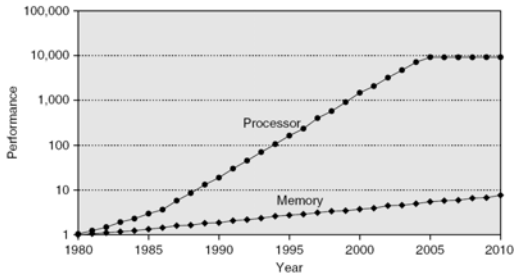
(a) Memory hierarchy for server

Component	Size	Speed
Register reference	500 bytes	500 ps
Level 1 Cache reference	64 KB	2 ns
Level 2 Cache reference	256 KB	10-20 ns
Memory reference	256-512 MB	50-100 ns
FLASH memory reference	4-8 GB	25-50 us

(b) Memory hierarchy for a personal mobile device

MK Copyright © 2012, Elsevier Inc. All rights reserved. 3

Memory Performance Gap



Memory Hierarchy Design

- Memory hierarchy design becomes more crucial with recent multi-core processors:
 - Aggregate peak bandwidth grows with # cores:
 - Intel Core i7 can generate two references per core per clock
 - Four cores and 3.2 GHz clock
 - 25.6 billion 64-bit data references/second +
 - 12.8 billion 128-bit instruction references
 - = 409.6 GB/s!
 - DRAM bandwidth is only 6% of this (25 GB/s)
 - Requires:
 - Multi-port, pipelined caches
 - Two levels of cache per core
 - Shared third-level cache on chip

Performance and Power

- High-end microprocessors have >10 MB on-chip cache
 - Consumes large amount of area and power budget

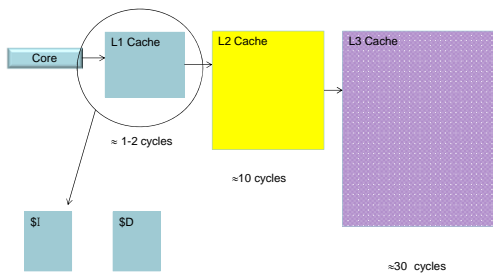
Terminology

- **A Block:** The smallest unit of information transferred between two levels.
- **Hit:** Item is found in some block in the upper level (example: Block X)
- **Miss:** Item needs to be retrieved from a block in the lower level (Block Y)
 - **Miss Rate** = 1 - (Hit Rate)
 - **Miss Penalty:** Time to replace a block in the upper level + Time to deliver the block the processor

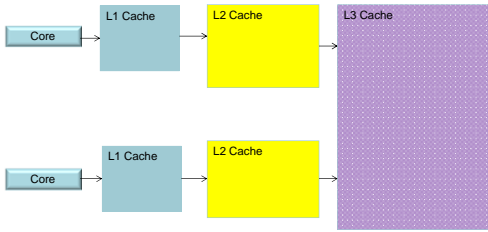
Cache operation

- Questions
 1. Where a block be placed in the cache (**placement**)
 2. How is a block is found if it is in the cache (**identification**)
 3. Which block should be replaced on a miss (**replacement**)
 4. What happens on a write (**write strategy**)

The Cache Hierarchy



The Cache Hierarchy



Cache Organization: Placement

- 1 **Direct mapped cache:** A block can be placed in only one location (cache block frame), given by the mapping function:
$$\text{index} = (\text{Block address}) \text{ MOD } (\text{Number of blocks in cache})$$
- 2 **Fully associative cache:** A block can be placed anywhere in cache. (no mapping function).
- 3 **Set associative cache:** A block can be placed in a restricted set of places, or cache block frames. A set is a group of block frames in the cache. A block is first mapped onto the set and then it can be placed anywhere within the set. The set in this case is chosen by:
$$\text{index} = (\text{Block address}) \text{ MOD } (\text{Number of sets in cache})$$

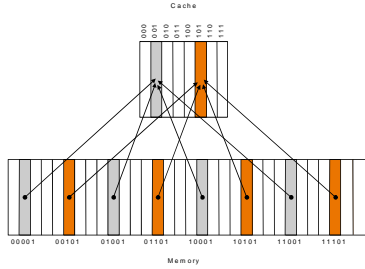
If there are n blocks in a set the cache placement is called n -way set-associative.

Cache Miss

- **Compulsory:** The very first access to a block is always a miss— Occurs even if you have an infinite cache
- **Capacity:** The cache is not big enough to hold all the blocks required for the execution of the program— A bigger cache helps
- **Conflict:** If not a fully associative, a block may be discarded and brought back again.

Cache Organization: Placement

Direct mapped Cache



Placement: DM

IK = 1024 Blocks
Each block = one word

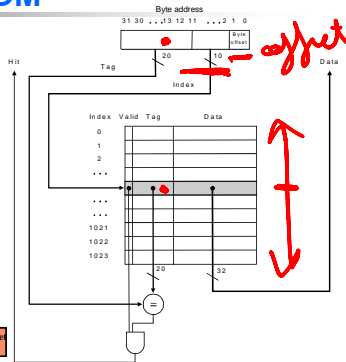
Can cache up to
 2^{32} bytes = 4 GB
of memory

Mapping function:

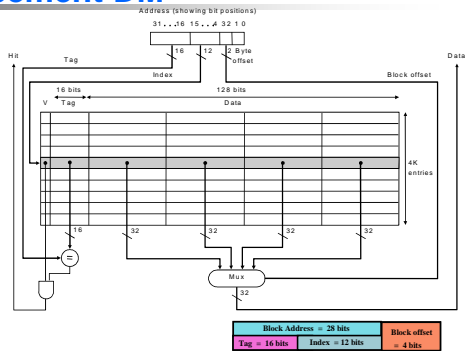
Cache Block frame number =
(Block address) MOD (1024)

i.e. index field or
10 low bit of block address

Block Address = 30 bits	Block offset = 2 bits
Tag = 20 bits	Index = 10 bits

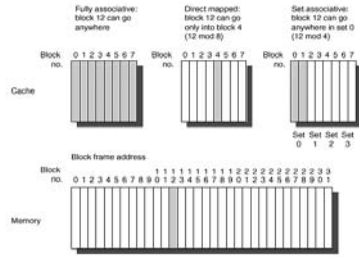


Placement DM



Block Address = 28 bits	Block offset = 4 bits
Tag = 16 bits	Index = 12 bits

Cache Organization

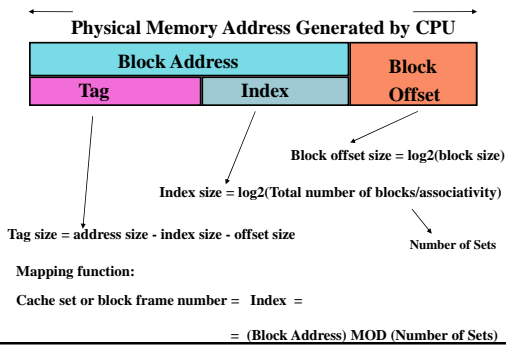


Cache Organization

- Each block frame in cache has an address tag.
- The tags of every cache block that might contain the required data are checked in parallel.
- A valid bit is added to the tag to indicate whether this entry contains a valid address.
- The address from the CPU to cache is divided into:
 - A block address, further divided into:
 - An index field to choose a block set in cache. (no index field when fully associative).
 - A tag field to search and match addresses in the selected set.
 - A block offset to select the data from the block.



Cache Organization



Cache Performance

- Assuming the following execution and cache parameters:
 - Cache miss penalty = 50 cycles
 - Normal instruction execution CPI ignoring memory stalls = 2.0 cycles
 - Miss rate = 2%
 - Average memory references/instruction = 1.33
- CPU time = $IC \times [CPI_{\text{execution}} + \text{Memory accesses/instruction} \times \text{Miss rate} \times \text{Miss penalty}] \times \text{Clock cycle time}$
- CPUtime with cache = $IC \times (2.0 + (1.33 \times 2\% \times 50)) \times \text{clock cycle time}$
- $= IC \times 3.33 \times \text{Clock cycle time}$
- Lower CPI execution increases the impact of cache miss clock cycles

Cache Performance

- 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 50 cycles.
- $CPI = CPI_{\text{execution}} + \text{mem stalls per instruction}$
- Mem Stalls per instruction =
 $\text{Mem accesses per instruction} \times \text{Miss rate} \times \text{Miss penalty}$
- Mem accesses per instruction = $1 + .3 = 1.3$
- Mem Stalls per instruction = $1.3 \times .015 \times 50 = 0.975$
- $CPI = 1.1 + .975 = 2.075$
- The ideal memory CPU with no misses is $2.075/1.1 = 1.88$ times faster

Cache Performance

- 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:
 - Miss penalty = $50 \times 2 = 100$ cycles.
 - $CPI = 1.1 + 1.3 \times .015 \times 100 = 1.1 + 1.95 = 3.05$
 - Speedup = $(CPI_{\text{old}} \times C_{\text{old}}) / (CPI_{\text{new}} \times C_{\text{new}})$
 $= 2.075 \times 2 / 3.05 = 1.36$
- 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.

Write Policy

Write Allocate:

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

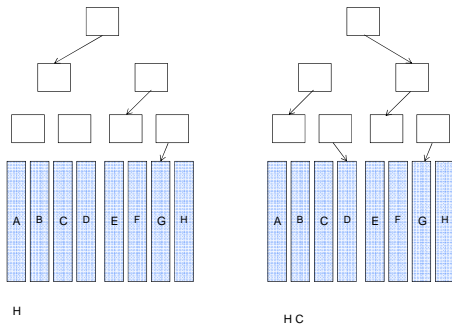
No-Write Allocate:

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

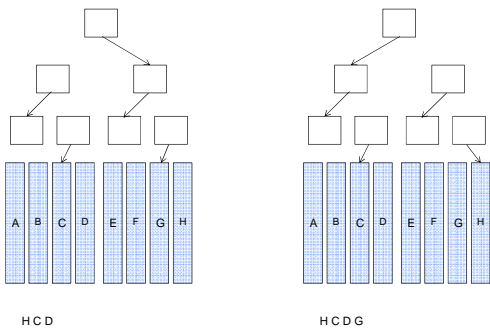
LRU

- A list to keep track of the order of access to every block in the set.
- The least recently used block is replaced (if needed).
- How many bits we need for that?

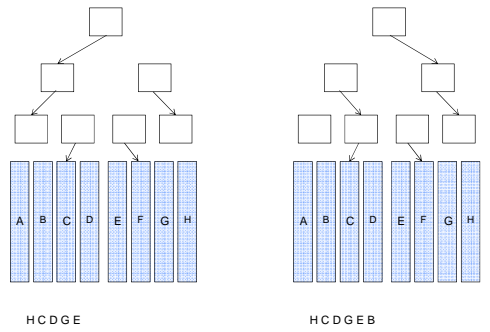
Pseudo LRU



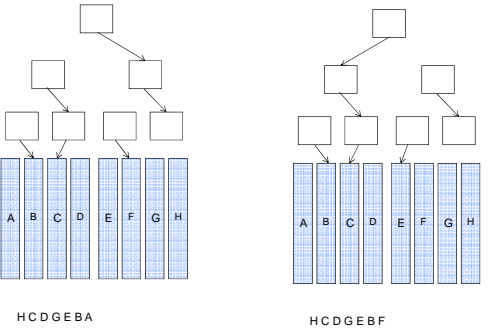
Pseudo LRU



Pseudo LRU



Pseudo LRU



Example

- Which has a lower miss rate 16KB cache for both instruction or data, or a combined 32KB cache? (0.64%, 6.47%, 1.99%).
- Assume hit=1cycle and miss =50 cycles. 75% of memory references are instruction fetch. *reads*
- Miss rate of split cache= $0.75 \times 0.64\% + 0.25 \times 6.47\% = 2.1\%$
- Slightly worse than 1.99% for combined cache. But, what about average memory access time?
- Split cache: $75\%(1+0.64\% \times 50) + 25\%(1+6.47\% \times 50) = 2.05$ cycles.
- Combined cache: $75\%(1+1.99\% \times 50) + 25\%(1+1.99\% \times 50) = 2.24$

Extra cycle for load/store

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.
 $CPI = CPI_{\text{execution}} + \text{mem stalls per instruction}$
% reads = $1.15/1.3 = 88.5\%$ % writes = $.15/1.3 = 11.5\%$

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.

1.3 mem ref/inst *50+50*

$$\frac{1.5}{1.02} (1.3 \times 50 \times 0.9 + 1.3 \times 0.1 \times 100)$$

Example

- CPU with $CPI_{\text{execution}} = 1.1$ running at clock rate = 500 MHz
- 1.3 memory accesses per instruction.
- L_1 cache operates at 500 MHz with a miss rate of 5%
- L_2 cache operates at 250 MHz with local miss rate 40%, ($T_2 = 2$ cycles)
- Memory access penalty, $M = 100$ cycles. Find CPI.

$$1.3 \left(\frac{5}{100} \times 0.6 \times 2 + \frac{5}{100} \times 0.4 \times 100 \right)$$

Example

- 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 with a miss rate of 1-H1 = 5%
 - Write through to L_2 with perfect write buffer with write allocate
- For L_2 :
 - Cache operates at 250 MHz with local miss rate 1-H2 = 40%, ($T_2 = 2$ cycles)
 - Write back to main memory with write allocate
 - Probability a cache block is dirty = 10%
- Memory access penalty, $M = 100$ cycles. Find CPI.

$$0.05 \left(0.6 \times 2 + 0.4 \times 0.9 \times 100 + 0.4 \times 0.1 \times 700 \right)$$

Example

- CPU with $CPI_{\text{execution}} = 1.1$ running at clock rate = 500 MHz
- 1.3 memory accesses per instruction.
- L_1 cache operates at 500 MHz with a miss rate of 5%
- L_2 cache operates at 250 MHz with a local miss rate 40%, ($T_2 = 2$ cycles)
- L_3 cache operates at 100 MHz with a local miss rate 50%, ($T_3 = 5$ cycles)
- Memory access penalty, $M = 100$ cycles. Find CPI.

HW
