
5DV008
Computer Architecture
Umeå University
Department of Computing Science

Stephen J. Hegner

Topic 5: The Memory Hierarchy
Part A: Caches

These slides are mostly taken verbatim, or with minor changes, from those prepared by

Mary Jane Irwin (www.cse.psu.edu/~mji)
of The Pennsylvania State University

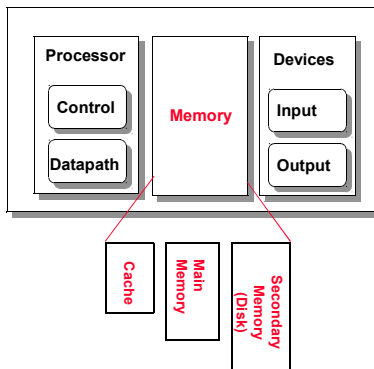
[Adapted from *Computer Organization and Design, 4th Edition*,
Patterson & Hennessy, © 2008, MK]

Key to the Slides

□ The source of each slide is coded in the footer on the right side:

- **Inwin CSE331** = slide by Mary Jane Irwin from the course CSE331 (Computer Organization and Design) at Pennsylvania State University.
- **Inwin CSE431** = slide by Mary Jane Irwin from the course CSE431 (Computer Architecture) at Pennsylvania State University.
- **Hegner UU** = slide by Stephen J. Hegner at Umeå University.

Review: Major Components of a Computer

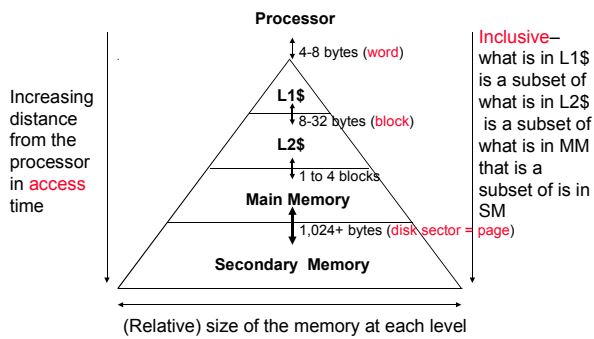


The Memory Hierarchy: Terminology

- **Block** (or line): the minimum unit of information that is present (or not) in a cache
- **Hit Rate**: the fraction of memory accesses found in a level of the memory hierarchy
 - **Hit Time**: Time to access that level which consists of
Time to access the block + Time to determine hit/miss
- **Miss Rate**: the fraction of memory accesses *not* found in a level of the memory hierarchy $\Rightarrow 1 - (\text{Hit Rate})$
 - **Miss Penalty**: Time to replace a block in that level with the corresponding block from a lower level which consists of
Time to access the block in the lower level + Time to transmit that block to the level that experienced the miss + Time to insert the block in that level + Time to pass the block to the requestor

Hit Time \ll Miss Penalty

Characteristics of the Memory Hierarchy



How is the Hierarchy Managed?

- registers \leftrightarrow memory
 - by compiler (programmer?)
- cache \leftrightarrow main memory
 - by the cache controller hardware
- main memory \leftrightarrow disks
 - by the operating system (virtual memory)
 - virtual to physical address mapping assisted by the hardware (TLB)
 - by the programmer (files)

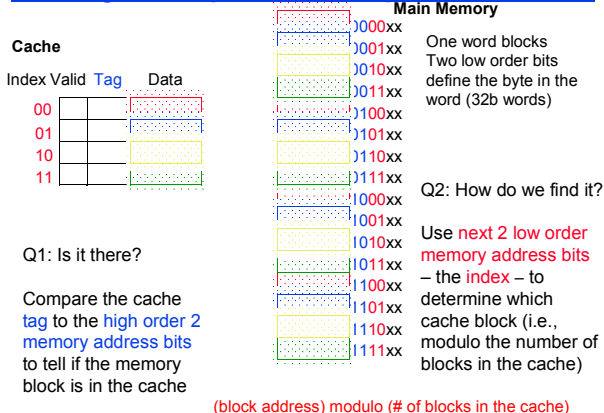
Cache Basics

- Two questions to answer (in hardware):
 - Q1: How do we know if a data item is in the cache?
 - Q2: If it is, how do we find it?
- Direct mapped
 - Each memory block is mapped to exactly one block in the cache
 - lots of lower level blocks must **share** blocks in the cache
 - Address mapping (to answer Q2):
 $(\text{block address}) \bmod (\# \text{ of blocks in the cache})$
 - Have a **tag** associated with each cache block that contains the address information (the upper portion of the address) required to identify the block (to answer Q1)

SDV008 20090212 1:5A sl:13

Irwin CSE431 PSU

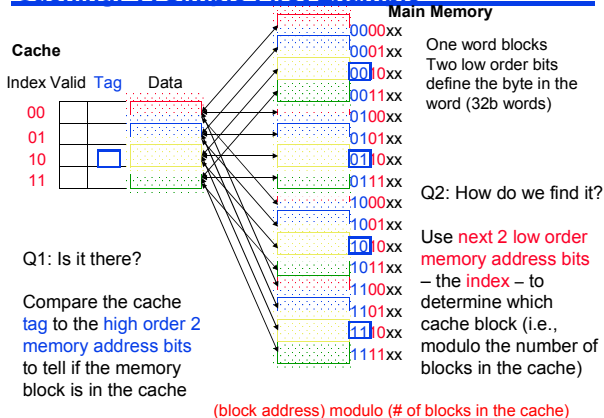
Caching: A Simple First Example



SDV008 20090212 1:5A sl:14

Irwin CSE431 PSU

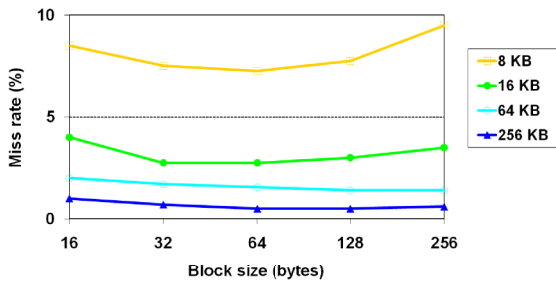
Caching: A Simple First Example



SDV008 20090212 1:5A sl:15

Irwin CSE431 PSU

Miss Rate vs Block Size vs Cache Size



- Miss rate goes up if the block size becomes a significant fraction of the cache size because the number of blocks that can be held in the same size cache is smaller (increasing **capacity** misses)

SDV008 20090212 1:5A sl:22

Irwin CSE431 PSU

Cache Field Sizes

- The number of bits in a cache includes both the storage for data and for the tags
 - 32-bit byte address
 - For a direct mapped cache with 2^n blocks, n bits are used for the index
 - For a block size of 2^m words (2^{m+2} bytes), m bits are used to address the word within the block and 2 bits are used to address the byte within the word
- What is the size of the tag field?
- The total number of bits in a direct-mapped cache is then $2^n \times (\text{block size} + \text{tag field size} + \text{valid field size})$
- How many total bits are required for a direct mapped cache with 16KB of data and 4-word blocks assuming a 32-bit address?

SDV008 20090212 1:5A sl:23

Irwin CSE431 PSU

Handling Cache Hits

- Read hits (I\$ and D\$)
 - this is what we want!
- Write hits (D\$ only)
 - require the cache and memory to be **consistent**
 - always write the data into both the cache block and the next level in the memory hierarchy (**write-through**)
 - writes run at the speed of the next level in the memory hierarchy – so slow! – or can use a **write buffer** and stall only if the write buffer is full
 - allow cache and memory to be **inconsistent**
 - write the data only into the cache block (**write-back** the cache block to the next level in the memory hierarchy when that cache block is “evicted”)
 - need a **dirty** bit for each data cache block to tell if it needs to be written back to memory when it is evicted – can use a **write buffer** to help “buffer” write-backs of dirty blocks

SDV008 20090212 1:5A sl:24

Irwin CSE431 PSU

Sources of Cache Misses

- **Compulsory** (cold start or process migration, first reference):
 - First access to a block, “cold” fact of life, not a whole lot you can do about it. If you are going to run “millions” of instruction, compulsory misses are insignificant
 - Solution: increase block size (increases miss penalty; very large blocks could increase miss rate)
- **Capacity:**
 - Cache cannot contain all blocks accessed by the program
 - Solution: increase cache size (may increase access time)
- **Conflict** (collision):
 - Multiple memory locations mapped to the same cache location
 - Solution 1: increase cache size
 - Solution 2: increase associativity (stay tuned) (may increase access time)

Handling Cache Misses (Single Word Blocks)

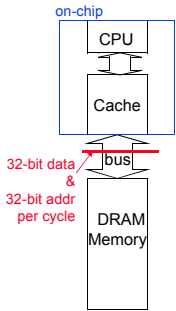
- **Read misses (I\$ and D\$)**
 - **stall** the pipeline, fetch the block from the next level in the memory hierarchy, install it in the cache and send the requested word to the processor, then let the pipeline resume
- **Write misses (D\$ only)**
 1. **Write allocate for multiple-word blocks** – **stall** the pipeline, fetch the block from next level in the memory hierarchy, install it in the cache (which may involve having to evict a dirty block if using a write-back cache), write the word from the processor to the cache, then let the pipeline resume, or
 2. **Faster write allocate for single-word blocks** – just write the word into the cache updating both the tag and data, no need to check for cache hit, no need to stall, or
 3. **No-write allocate** – skip the cache write (but must invalidate that cache block since it will now hold stale data) and just write the word to the write buffer (and eventually to the next memory level), no need to stall if the write buffer isn't full

Multiword Block Considerations

- **Read misses (I\$ and D\$)**
 - Processed the same as for single word blocks – a miss returns the entire block from memory
 - Miss penalty grows as block size grows
 - **Early restart** – processor resumes execution as soon as the requested word of the block is returned
 - **Requested word first** – requested word is transferred from the memory to the cache (and processor) first
 - **Nonblocking cache** – allows the processor to continue to access the cache while the cache is handling an earlier miss
- **Write misses (D\$)**
 - If using write allocate must *first* fetch the block from memory and then write the word to the block (or could end up with a “garbled” block in the cache (e.g., for 4 word blocks, a new tag, one word of data from the new block, and three words of data from the old block))

Memory Systems that Support Caches

- The off-chip interconnect and memory architecture can affect overall system performance in dramatic ways



One-word-wide organization (one-word-wide bus and one-word-wide memory)

- Assume
 1. 1 memory bus clock cycle to send the addr
 2. 15 memory bus clock cycles to get the 1st word in the block from DRAM (row cycle time), 5 memory bus clock cycles for 2nd, 3rd, 4th words (column access time)
 3. 1 memory bus clock cycle to return a word of data
- Memory-Bus to Cache bandwidth
 - number of bytes accessed from memory and transferred to cache/CPU per memory bus clock cycle

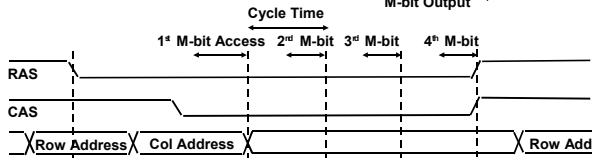
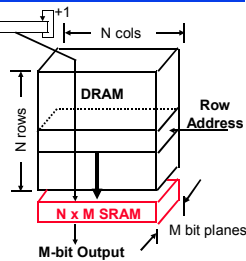
SDV008 20090212 1:5A sl:28

Irwin CSE431 PSU

Review: (DDR) SDRAM Operation

- After a row is read into the SRAM register

- Input CAS as the starting "burst" address along with a burst length
- Transfers a burst of data (ideally a cache block) from a series of sequential addr's within that row
 - The memory bus clock controls transfer of successive words in the burst



SDV008 20090212 1:5A sl:29

Irwin CSE431 PSU

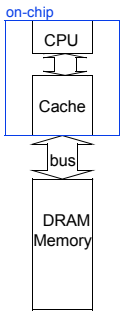
DRAM Size Increase

- Add a table like figure 5.12 to show DRAM growth since 1980

SDV008 20090212 1:5A sl:30

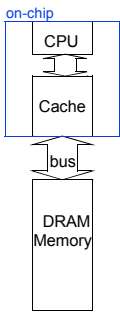
Irwin CSE431 PSU

One Word Wide Bus. One Word Blocks



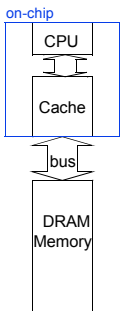
- If the block size is one word, then for a memory access due to a cache miss, the pipeline will have to stall for the number of cycles required to return one data word from memory
 - cycle to send address
 - cycles to read DRAM
 - cycle to return data
 - total clock cycles miss penalty
- Number of bytes transferred per clock cycle (bandwidth) for a single miss is
 - bytes per memory bus clock cycle

One Word Wide Bus. One Word Blocks



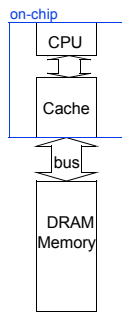
- If the block size is one word, then for a memory access due to a cache miss, the pipeline will have to stall for the number of cycles required to return one data word from memory
 - 1 memory bus clock cycle to send address
 - 15 memory bus clock cycles to read DRAM
 - 1 memory bus clock cycle to return data
 - 17 total clock cycles miss penalty
- Number of bytes transferred per clock cycle (bandwidth) for a single miss is
 - $4/17 = 0.235$ bytes per memory bus clock cycle

One Word Wide Bus. Four Word Blocks



- What if the block size is four words and each word is in a different DRAM row?
 - cycle to send 1st address
 - cycles to read DRAM
 - cycles to return last data word
 - total clock cycles miss penalty
- Number of bytes transferred per clock cycle (bandwidth) for a single miss is
 - bytes per clock

One Word Wide Bus. Four Word Blocks

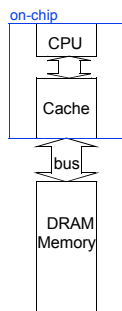


- What if the block size is four words and each word is in a different DRAM row?
 - 1 cycle to send 1st address
 - $4 \times 15 = 60$ cycles to read DRAM
 - 1 cycles to return last data word
 - 62 total clock cycles miss penalty



- Number of bytes transferred per clock cycle (bandwidth) for a single miss is
 - $(4 \times 4) / 62 = 0.258$ bytes per clock

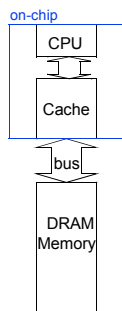
One Word Wide Bus. Four Word Blocks



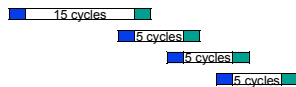
- What if the block size is four words and all words are in the same DRAM row?
 - 1 cycle to send 1st address
 - 3 cycles to read DRAM
 - 1 cycles to return last data word
 - 5 total clock cycles miss penalty

- Number of bytes transferred per clock cycle (bandwidth) for a single miss is
 - 4 bytes per clock

One Word Wide Bus. Four Word Blocks

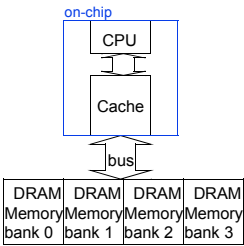


- What if the block size is four words and all words are in the same DRAM row?
 - 1 cycle to send 1st address
 - $15 + 3 \times 5 = 30$ cycles to read DRAM
 - 1 cycles to return last data word
 - 32 total clock cycles miss penalty



- Number of bytes transferred per clock cycle (bandwidth) for a single miss is
 - $(4 \times 4) / 32 = 0.5$ bytes per clock

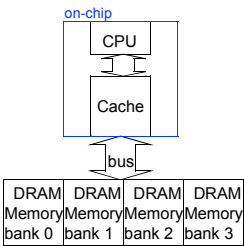
Interleaved Memory. One Word Wide Bus



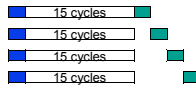
- For a block size of four words
 - _____ cycle to send 1st address
 - _____ cycles to read DRAM banks
 - _____ cycles to return last data word
 - _____ total clock cycles miss penalty

- Number of bytes transferred per clock cycle (bandwidth) for a single miss is _____ bytes per clock

Interleaved Memory. One Word Wide Bus



- For a block size of four words
 - 1 cycle to send 1st address
 - 15 cycles to read DRAM banks
 - 4*1 = 4 cycles to return last data word
 - 20 total clock cycles miss penalty



- Number of bytes transferred per clock cycle (bandwidth) for a single miss is _____ bytes per clock
- (4 x 4)/20 = 0.8

DRAM Memory System Summary

- It is important to match the cache characteristics
 - caches access one block at a time (usually more than one word)
- with the DRAM characteristics
 - use DRAMs that support fast multiple word accesses, preferably ones that match the block size of the cache
- with the memory-bus characteristics
 - make sure the memory-bus can support the DRAM access rates and patterns
 - with the goal of increasing the Memory-Bus to Cache bandwidth

Measuring Cache Performance

- Assuming cache hit costs are included as part of the normal CPU execution cycle, then

$$\begin{aligned} \text{CPU time} &= \text{IC} \times \text{CPI} \times \text{CC} \\ &= \text{IC} \times (\underbrace{\text{CPI}_{\text{ideal}} + \text{Memory-stall cycles}}_{\text{CPI}_{\text{stall}}}) \times \text{CC} \end{aligned}$$

- Memory-stall cycles come from cache misses (a sum of read-stalls and write-stalls)

$$\text{Read-stall cycles} = \text{reads/program} \times \text{read miss rate} \times \text{read miss penalty}$$

$$\begin{aligned} \text{Write-stall cycles} &= (\text{writes/program} \times \text{write miss rate} \times \text{write miss penalty}) \\ &\quad + \text{write buffer stalls} \end{aligned}$$

- For write-through caches, we can simplify this to

$$\text{Memory-stall cycles} = \text{accesses/program} \times \text{miss rate} \times \text{miss penalty}$$

Impacts of Cache Performance

- Relative cache penalty increases as processor performance improves (faster clock rate and/or lower CPI)

- The memory speed is unlikely to improve as fast as processor cycle time. When calculating $\text{CPI}_{\text{stall}}$, the cache miss penalty is measured in processor clock cycles needed to handle a miss
- The lower the $\text{CPI}_{\text{ideal}}$, the more pronounced the impact of stalls

- A processor with a $\text{CPI}_{\text{ideal}}$ of 2, a 100 cycle miss penalty, 36% load/store instr's, and 2% I\$ and 4% D\$ miss rates
Memory-stall cycles = $2\% \times 100 + 36\% \times 4\% \times 100 = 3.44$

$$\text{So } \text{CPI}_{\text{stalls}} = 2 + 3.44 = \mathbf{5.44}$$

more than twice the $\text{CPI}_{\text{ideal}}$!

- What if the $\text{CPI}_{\text{ideal}}$ is reduced to 1? 0.5? 0.25?
- What if the D\$ miss rate went up 1%? 2%?
- What if the processor clock rate is doubled (doubling the miss penalty)?

Average Memory Access Time (AMAT)

- A larger cache will have a longer access time. An increase in hit time will likely add another stage to the pipeline. At some point the increase in hit time for a larger cache will overcome the improvement in hit rate leading to a decrease in performance.
- Average Memory Access Time (AMAT) is the average to access memory considering both hits and misses

$$\text{AMAT} = \text{Time for a hit} + \text{Miss rate} \times \text{Miss penalty}$$

- What is the AMAT for a processor with a 20 psec clock, a miss penalty of 50 clock cycles, a miss rate of 0.02 misses per instruction and a cache access time of 1 clock cycle?

Reducing Cache Miss Rates #1

- Allow more flexible block placement
 - In a **direct mapped cache** a memory block maps to exactly one cache block
 - At the other extreme, could allow a memory block to be mapped to *any* cache block – **fully associative cache**
 - A compromise is to divide the cache into **sets** each of which consists of n “ways” (**n-way set associative**). A memory block maps to a unique set (specified by the index field) and can be placed in any way of that set (so there are n choices)

(block address) modulo (# sets in the cache)

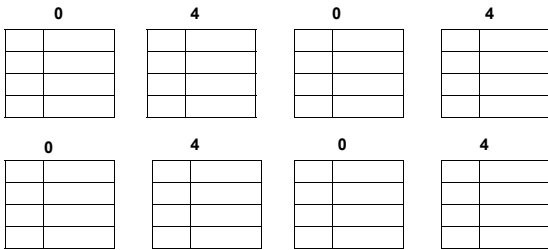
SDV008 20090212 1:5A sl:43

Irwin CSE431 PSU

Another Reference String Mapping

- Consider the main memory word reference string

Start with an empty cache - all blocks initially marked as not valid 0 4 0 4 0 4 0 4



SDV008 20090212 1:5A sl:44

Irwin CSE431 PSU

Another Reference String Mapping

- Consider the main memory word reference string

Start with an empty cache - all blocks initially marked as not valid 0 4 0 4 0 4 0 4



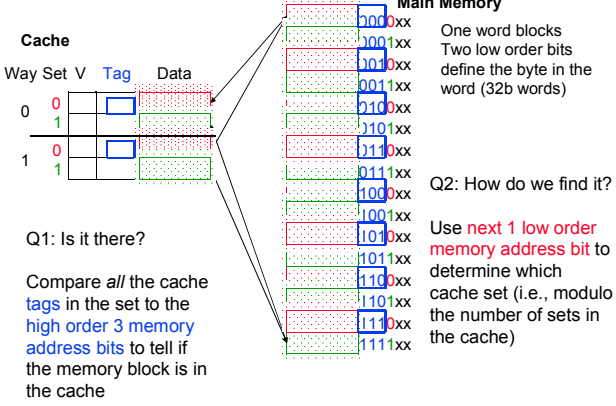
• 8 requests, 8 misses

- Ping pong effect due to **conflict** misses - two memory locations that map into the same cache block

SDV008 20090212 1:5A sl:45

Irwin CSE431 PSU

Set Associative Cache Example

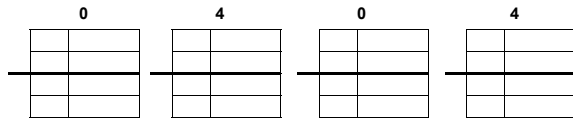


SDV008 20090212 1:5A sl:46

Irwin CSE431 PSU

Another Reference String Mapping

□ Consider the main memory word reference string
Start with an empty cache - all blocks initially marked as not valid

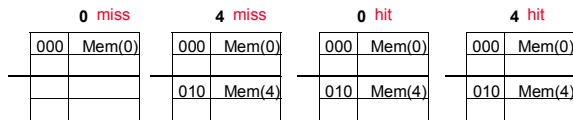


SDV008 20090212 1:5A sl:47

Irwin CSE431 PSU

Another Reference String Mapping

□ Consider the main memory word reference string
Start with an empty cache - all blocks initially marked as not valid



● 8 requests, 2 misses

□ Solves the ping pong effect in a direct mapped cache due to **conflict** misses since now two memory locations that map into the same cache set can co-exist!

SDV008 20090212 1:5A sl:48

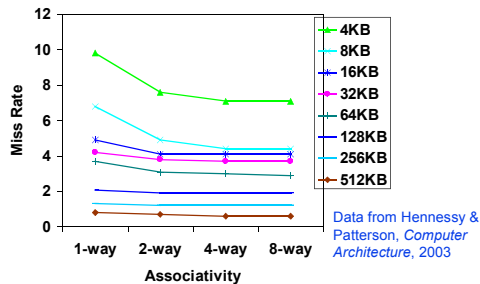
Irwin CSE431 PSU

Costs of Set Associative Caches

- When a miss occurs, which way's block do we pick for replacement?
 - Least Recently Used (LRU): the block replaced is the one that has been unused for the longest time
 - Must have hardware to keep track of when each way's block was used relative to the other blocks in the set
 - For 2-way set associative, takes **one bit per set** → set the bit when a block is referenced (and reset the other way's bit)
- N-way set associative cache costs
 - N comparators (delay and area)
 - MUX delay (set selection) before data is available
 - Data available **after** set selection (and Hit/Miss decision). In a direct mapped cache, the cache block is available **before** the Hit/Miss decision
 - So its not possible to just assume a hit and continue and recover later if it was a miss

Benefits of Set Associative Caches

- The choice of direct mapped or set associative depends on the cost of a miss versus the cost of implementation



- Largest gains are in going from direct mapped to 2-way (20%+ reduction in miss rate)

Reducing Cache Miss Rates #2

1. Use multiple levels of caches
 - With advancing technology have more than enough room on the die for bigger L1 caches or for a second level of caches – normally a **unified** L2 cache (i.e., it holds both instructions and data) and in some cases even a unified L3 cache
 - For our example, CPI_{ideal} of 2, 100 cycle miss penalty (to main memory) and a 25 cycle miss penalty (to UL2\$), 36% load/stores, a 2% (4%) L1 I\$ (D\$) miss rate, add a 0.5% UL2\$ miss rate

$$CPI_{stalls} = 2 + .02 \times 25 + .36 \times .04 \times 25 + .005 \times 100 + .36 \times .005 \times 100 = 3.54$$

(as compared to 5.44 with no L2\$)

Multilevel Cache Design Considerations

- Design considerations for L1 and L2 caches are very different
 - Primary cache should focus on **minimizing hit time** in support of a shorter clock cycle
 - Smaller with smaller block sizes
 - Secondary cache(s) should focus on **reducing miss rate** to reduce the penalty of long main memory access times
 - Larger with larger block sizes
 - Higher levels of associativity
- The miss penalty of the L1 cache is significantly reduced by the presence of an L2 cache – so it can be smaller (i.e., faster) but have a higher miss rate
- For the L2 cache, hit time is less important than miss rate
 - The L2\$ hit time determines L1's miss penalty
 - L2\$ local miss rate >> than the global miss rate

Using the Memory Hierarchy Well

- Include plots from Figure 5.18

Two Machines' Cache Parameters

| | Intel Nehalem | AMD Barcelona |
|------------------------------|---|---|
| L1 cache organization & size | Split I\$ and D\$; 32KB for each per core; 64B blocks | Split I\$ and D\$; 64KB for each per core; 64B blocks |
| L1 associativity | 4-way (I), 8-way (D) set assoc.; ~LRU replacement | 2-way set assoc.; LRU replacement |
| L1 write policy | write-back, write-allocate | write-back, write-allocate |
| L2 cache organization & size | Unified; 256KB (0.25MB) per core; 64B blocks | Unified; 512KB (0.5MB) per core; 64B blocks |
| L2 associativity | 8-way set assoc.; ~LRU | 16-way set assoc.; ~LRU |
| L2 write policy | write-back | write-back |
| L2 write policy | write-back, write-allocate | write-back, write-allocate |
| L3 cache organization & size | Unified; 8192KB (8MB) shared by cores; 64B blocks | Unified; 2048KB (2MB) shared by cores; 64B blocks |
| L3 associativity | 16-way set assoc. | 32-way set assoc.; evict block shared by fewest cores |
| L3 write policy | write-back, write-allocate | write-back; write-allocate |

Two Older Machines' Cache Parameters

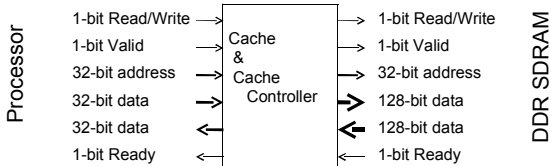
| | Intel P4 | AMD Opteron |
|------------------|--|------------------------------|
| L1 organization | Split I\$ and D\$ | Split I\$ and D\$ |
| L1 cache size | 8KB for D\$, 96KB for trace cache (~I\$) | 64KB for each of I\$ and D\$ |
| L1 block size | 64 bytes | 64 bytes |
| L1 associativity | 4-way set assoc. | 2-way set assoc. |
| L1 replacement | ~LRU | LRU |
| L1 write policy | write-through | write-back |
| L2 organization | Unified | Unified |
| L2 cache size | 512KB | 1024KB (1MB) |
| L2 block size | 128 bytes | 64 bytes |
| L2 associativity | 8-way set assoc. | 16-way set assoc. |
| L2 replacement | ~LRU | ~LRU |
| L2 write policy | write-back | write-back |

SDV008 20090212 1:5A sl:58

Irwin CSE431 PSU

FSM Cache Controller

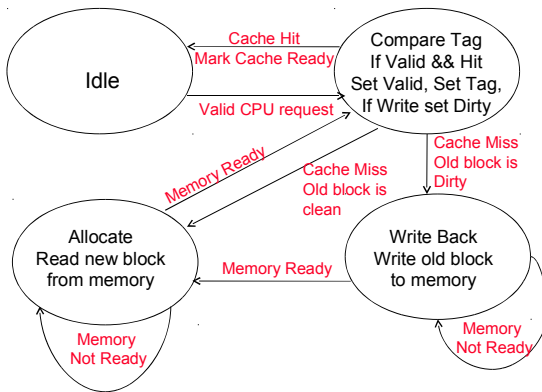
- Key characteristics for a simple L1 cache
 - Direct mapped
 - Write-back using write-allocate
 - Block size of 4 32-bit words (so 16B); Cache size of 16KB (so 1024 blocks)
 - 18-bit tags, 10-bit index, 2-bit block offset, 2-bit byte offset, dirty bit, valid bit, LRU bits (if set associative)



SDV008 20090212 1:5A sl:59

Irwin CSE431 PSU

Four State Cache Controller



SDV008 20090212 1:5A sl:60

Irwin CSE431 PSU

Summary: Improving Cache Performance

0. Reduce the time to hit in the cache

- smaller cache
- direct mapped cache
- smaller blocks
- for writes
 - no write allocate – no "hit" on cache, just write to write buffer
 - write allocate – to avoid two cycles (first check for hit, then write)
pipeline writes via a delayed write buffer to cache

1. Reduce the miss rate

- bigger cache
- more flexible placement (increase associativity)
- larger blocks (16 to 64 bytes typical)
- victim cache – small buffer holding most recently discarded blocks

Summary: Improving Cache Performance

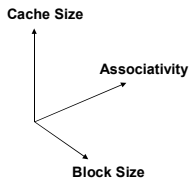
2. Reduce the miss penalty

- smaller blocks
- use a write buffer to hold dirty blocks being replaced so don't have to wait for the write to complete before reading
- check write buffer (and/or victim cache) on read miss – may get lucky
- for large blocks fetch critical word first
- use multiple cache levels – L2 cache not tied to CPU clock rate
- faster backing store/improved memory bandwidth
 - wider buses
 - memory interleaving, DDR SDRAMs

Summary: The Cache Design Space

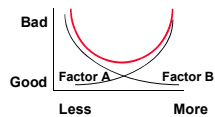
□ Several interacting dimensions

- cache size
- block size
- associativity
- replacement policy
- write-through vs write-back
- write allocation



□ The optimal choice is a compromise

- depends on access characteristics
 - workload
 - use (I-cache, D-cache, TLB)
- depends on technology / cost



□ Simplicity often wins
