Memory Hierarchy Design - Cache Memory

STUDENTS-HUB.com

Presentation Outline

- Memory Hierarchy and the need for Cache Memory
- The Basics of Caches
- Multilevel Caches
- Cache Performance and Memory Stall Cycles

Processor-Memory Performance Gap



- ✤ 1980 No cache in microprocessor
- ✤ 1995 Two-level cache on microprocessor

STUDENTS-HUB.com

Motivation for Memory Hierarchy

Programmers want unlimited amounts of memory with low latency

♦ But fast memory more expensive than slower memory

Solution: small fast memory + big slow memory

= Looks like a big fast memory



The Need for Cache Memory

Widening speed gap between CPU and main memory

- ♦ Processor operation takes less than 1 ns
- \diamond Main memory requires more than 50 ns to access
- Each instruction involves at least one memory access
 - \diamond One memory access to fetch the instruction
 - ♦ A second memory access for load and store instructions
- Memory bandwidth limits the instruction execution rate
- Cache memory can help bridge the CPU-memory gap
- Cache memory is small in size but fast

STUDENTS-HUB.com

Model of Memory Hierarchy



STUDENTS-HUB.com

Typical Memory Hierarchy

- Registers are at the top of the hierarchy
 - ♦ Typical size < 1 KB</p>
 - ♦ Access time < 0.5 ns</p>
- ✤ Level 1 Cache (8 64 KiB)
 - ♦ Access time: 1 ns
- ✤ L2 Cache (1 MiB 8 MiB)
 - ♦ Access time: 3 10 ns
- ✤ Main Memory (8 32 GiB)
 - ♦ Access time: 40 50 ns
- Disk Storage (> 200 GB)
 - \diamond Access time: 5 10 ms



STUDENTS-HUB.com

Typical Memory Hierarchy



STUDENTS-HUB.com

Principle of Locality of Reference

- Programs access small portion of their address space
 - \diamond At any time, only a small set of instructions & data is needed

Temporal Locality (in time)

- ♦ If an item is accessed, probably it will be accessed again soon
- ♦ Same loop instructions are fetched each iteration
- ♦ Same procedure may be called and executed many times

Spatial Locality (in space)

- ♦ Tendency to access contiguous instructions/data in memory
- ♦ Sequential execution of Instructions
- ♦ Traversing arrays element by element

Memory Hierarchy Operation

- If an instruction or operand is required by the CPU, the levels of the memory hierarchy are searched for the item starting with the level closest to the CPU (Level 1 cache): L1 Cache
 - ♦ If the item is found, it's delivered to the CPU resulting in <u>a cache hit</u> without searching lower levels.

 Hit rate for level one cache = H_1
 - ♦ If the item is missing from an upper level, resulting in <u>a cache miss</u>, the level just below is searched.

 Miss rate for level one cache = 1 Hit rate = 1 Hit
 - ♦ For systems with several levels of cache, the search continues with cache level 2, 3 etc.
 - ♦ If all levels of cache report a miss then main memory is accessed for the item.
 - CPU \leftrightarrow cache \leftrightarrow memory: <u>Managed by hardware</u>.
 - If the item is not found in main memory resulting in a page fault, then disk (virtual memory), is accessed for the item.
 - Memory ↔ disk: <u>Managed by the operating system</u> with hardware support

STUDENTS-HUB.com

Memory Hierarchy: 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) *
 - **Hit Rate:** The fraction of memory access found in the upper level. \diamond
 - Hit Time: Time to access the upper level which consists of (S)RAM access time + Time to determine hit/miss

Ideally = 1 Cycle

- Miss: Item needs to be retrieved from a block in the lower level (Block Y) *
 - **Miss Rate = 1 (Hit Rate)** | Miss rate for level one cache = $1 \text{Hit rate} = 1 \text{H}_1$ \diamond
 - Miss Penalty: Time to replace a block in the upper level + M \diamond

Time to deliver the missed block to the processor

Hit Time << Miss Penalty * Μ



Basic Cache Concepts

- Cache is the first level of the memory hierarchy once the address leaves the CPU and is searched first for the requested data.
- If the data requested by the CPU is present in the cache, it is retrieved from cache and the data access is <u>a cache hit</u> otherwise <u>a cache miss</u> and data must be read from main memory.
- On a cache miss a block of data must be brought in from main memory to cache to possibly <u>replace</u> an existing cache block.
- The allowed block addresses where blocks can be mapped (placed) into cache from main memory is determined by <u>cache placement strategy</u>.
- Locating a block of data in cache is handled by cache <u>block identification</u> <u>mechanism</u>: Tag matching.
- On a cache miss choosing the cache block being removed (replaced) is handled by the <u>block replacement strategy</u> in place.
- When a write to cache is requested, a number of main memory update
 STUDENTS FAGES exist as part of <u>the cache write policy</u>.
 Uploaded By: Jibreel Bornat

What is a Cache Memory?

- Small and fast (SRAM) memory technology
 - ♦ Stores the subset of instructions & data currently being accessed
- Used to reduce average access time to memory
- Caches exploit temporal locality by ...
 - ♦ Keeping recently accessed data closer to the processor
- Caches exploit spatial locality by ...
 - ♦ Moving blocks consisting of multiple contiguous words
- Goal is to achieve
 - Fast speed of cache memory access
 - ♦ Balance the cost of the memory system

STUDENTS-HUB.com

Cache Memories in the Datapath



STUDENTS-HUB.com

Almost Everything is a Cache!

- In computer architecture, almost everything is a cache!
- Registers: a cache on variables software managed
- First-level cache: a cache on second-level cache
- Second-level cache: a **cache on memory** (or L3 cache)
- Memory: a cache on hard disk
 - ♦ Stores recent programs and their data
 - ♦ Hard disk can be viewed as an extension to main memory
- Branch target and prediction buffer

Cache on branch target and prediction information

STUDENTS-HUB.com

Presentation Outline

Memory Hierarchy and the need for Cache Memory

The Basics of Caches

Multilevel Caches

Cache Performance and Memory Stall Cycles

Four Basic Questions on Caches

✤ Q1: Where can a block be placed in a cache?

♦ Block placement

- ♦ Direct Mapped, Set Associative, Fully Associative
- ✤ Q2: How is a block found in a cache?
 - ♦ Block identification
 - ♦ Block address, tag, index
- ✤ Q3: Which block should be replaced on a cache miss?
 - ♦ Block replacement
 - ♦ FIFO, Random, LRU
- ✤ Q4: What happens on a write?

♦ Write strategy

Write Back or Write Through cache (with Write Buffer) STUDENTS-HUB.com

Inside a Cache Memory



Cache Block (or Cache Line)

- ♦ Unit of data transfer between main memory and a cache
- ♦ Large block size → Less tag overhead + Burst transfer from DRAM
- \diamond Typically, cache block size = 64 bytes in recent caches

STUDENTS-HUB.com

Block Placement: Direct Mapped

Block: unit of data transfer between cache and memory

Direct Mapped Cache:

♦ A block can be placed in exactly one location in the cache



STUDENTS-HUB.com

Direct Mapped Cache



STUDENTS-HUB.com

64KB Direct Mapped Cache Example



Direct-Mapped Cache

- ✤ A memory address is divided into
 - Block address: identifies block in memory
 - ♦ Block offset: to access bytes within a block
- A block address is further divided into
 - Index: used for direct cache access
 - Tag: most-significant bits of block address

Index = Block Address **mod** Cache Blocks

- Tag must be stored also inside cache
 - ♦ For block identification

A valid bit is also required to indicate

♦ Whether a cache block is valid or not STUDENTS-HUB.com



Direct Mapped Cache - cont'd

Cache hit: block is stored inside cache

- ♦ Index is used to access cache block
- ♦ Address tag is compared against stored tag
- $\diamond\,$ If equal and cache block is valid then hit
- ♦ Otherwise: cache miss
- If number of cache blocks is 2^n
 - \diamond *n* bits are used for the cache index
- * If number of bytes in a block is 2^b
 - \diamond *b* bits are used for the block offset
- ✤ If 32 bits are used for an address
 - \Rightarrow 32 *n b* bits are used for the tag
- Cache data size = 2^{n+b} bytes STUDENTS-HUB.com



Mapping an Address to a Cache Block

Example

- ♦ Consider a direct-mapped cache with 256 blocks
- \diamond Block size = 16 bytes
- ♦ Compute tag, index, and byte offset of address: 0x01FFF8AC

Solution

♦ 32-bit address is divided into:



- 4-bit byte offset field, because block size = 2⁴ = 16 bytes
- 8-bit cache index, because there are $2^8 = 256$ blocks in cache
- 20-bit tag field
- \Rightarrow Byte offset = 0xC = 12 (least significant 4 bits of address)
- \diamond Cache index = 0x8A = 138 (next lower 8 bits of address)
- \Rightarrow Tag = 0x01FFF (upper 20 bits of address)

STUDENTS-HUB.com

Example on Cache Placement & Misses

- Consider a small direct-mapped cache with 32 blocks
 - \diamond Cache is initially empty, Block size = 16 bytes
 - \diamond The following memory addresses (in decimal) are referenced: 1000, 1004, 1008, 2548, 2552, 2556.
 - \diamond Map addresses to cache blocks and indicate whether hit or miss

Solution:

		23	5	4	
Solution:		Тад	Index	offset	
♦ 1000 = 0x3E8	cache ir	ndex =	0x1E		Miss (first access)
♦ 1004 = 0x3EC	cache ir	ndex =	0x1E		Hit
♦ 1008 = 0x3F0	cache ir	ndex =	0x1F		Miss (first access)
♦ 2548 = 0x9F4	cache ir	ndex =	0x1F		Miss (different tag)
♦ 2552 = 0x9F8	cache ir	ndex =	0x1F		Hit
♦ 2556 = 0x9FC	cache ir	ndex =	0x1F		Hit

STUDENTS-HUB.com

Fully Associative Cache

- A block can be placed anywhere in cache \Rightarrow no indexing
- ✤ If *m* blocks exist then
 - ♦ m comparators are needed to match tag



STUDENTS-HUB.com

Set-Associative Cache

- * A set is a group of blocks that can be indexed
- ✤ One set = *m* blocks → *m*-way set associative
- Set index field is *k*-bit long \rightarrow Number of Sets = 2^k
- Set index is decoded and only one set is examined

 - ♦ If address tag matches a stored tag within set then Cache Hit
 - ♦ Otherwise: Cache Miss
- ✤ Cache data size = $m \times 2^{k+b}$ bytes (2^b bytes per block)
- A direct-mapped cache has one block per set (m = 1)

A fully-associative cache has only one set (k = 0)
STUDENTS-HUB.com

4K Four-Way Set Associative Cache: MIPS Implementation Example



Handling a Cache Read



Replacement Policy

- Which block to replace on a cache miss?
- No choice for direct-mapped cache
- m choices for m-way set associative cache

Random replacement

- ♦ Candidate block is selected randomly
- ♦ One counter for all sets (0 to m 1): incremented on every cycle
- ♦ On a cache miss, replace block specified by counter

First In First Out (FIFO) replacement

- ♦ Replace oldest block in set (Round-Robin)
- ♦ One counter per set (0 to m 1): specifies oldest block to replace
- ♦ Counter is incremented on a cache miss

STUDENTS-HUB.com

LRU Replacement Policy

LRU: Replace Least Recently Used Block in a Set

- ♦ LRU state must be updated on every cache hit
- ♦ If m = 2, there are only 2 permutations → a single LRU bit is needed
- ♦ If m = 4 then 4! = 24 permutations. If m = 8 then 8! = 40320 permutations
- \diamond Pure LRU is difficult to implement for large *m*

Pseudo LRU Tree

♦ LRU approximation that requires (m - 1) LRU bits per set



Access Stream: B, C, A, D, E

D



After access to A 0 0 С В Α D

After access to E. E replaces B, Bits flipped along the path







Comparing Random, FIFO, and LRU

- Data cache misses per 1000 instructions
 - ♦ 10 SPEC2000 benchmarks on Alpha processor
 - \diamond Block size = 64 bytes
 - ♦ LRU outperforms FIFO and Random for a small cache
 - ♦ Little difference between LRU and Random for a large cache
- Random is the simplest to implement: one counter for all sets
- ✤ Pseudo LRU Tree requires (m-1) replacement bits per set

	2-way		4-way			8-way			
Size	LRU	Rand	FIFO	LRU	Rand	FIFO	LRU	Rand	FIFO
16 KB	114.1	117.3	115.5	111.7	115.1	113.3	109.0	111.8	110.4
64 KB	103.4	104.3	103.9	102.4	102.3	103.1	99.7	100.5	100.3
256 KB	92.2	92.1	92.5	92.1	92.1	92.5	92.1	92.1	92.5

STUDENTS-HUB.com

Write Policy

Write Through

- ♦ Writes update data cache and lower-level (L2) cache
- ♦ Cache control bit: only a Valid bit is needed
- ♦ Lower-level cache always has latest data, which simplifies data coherency
- \diamond Can always discard cached data when a block is replaced

Write Back

- ♦ Writes update data cache only
- ♦ Two cache control bits: Valid and Modified bits are required
- \diamond Modified bit indicates that cache block \neq memory block
- Modified cached block is written back when replaced
- ♦ Multiple writes to same cache block require only one write-back

♦ Uses less memory bandwidth than write-through and less energy STUDENTS-HUB.com
Uploaded By: Jibreel Bornat

Write Miss Policy

What happens on a write miss?

Write Allocate

- ♦ Allocate new block in data cache
- ♦ Write miss acts like a read miss, block is fetched then updated

Write No-Allocate

- \diamond Send data to lower-level (L2) cache
- \diamond Data cache is not modified
- Either write-miss policy can be combined with either write policy
- Write back caches typically use write allocate on a write miss
 - ♦ Reasoning: subsequent writes will be captured in the cache

Write-through caches might choose write no-allocate

Reasoning: writes must still go to lower level memory STUDENTS-HUB.com
Uploaded By: Jibreel Bornat

Write Buffer

Write buffer is a queue that holds: address + write-data (wdata)

- ♦ Write-through: all writes are sent to lower-level cache
- Buffer decouples the write from the memory bus writing
 - ♦ Write occurs without stalling processor, until buffer is full

Problem: write buffer may hold data on a read miss

- ♦ If address is found, return data value in write buffer
- ♦ Transfer block from lower level cache and update D-cache



STUDENTS-HUB.com
Write Through with Write No-Allocate



STUDENTS-HUB.com

Victim Buffer

- Used by modified evicted blocks and their addresses
 - \diamond When a modified block is replaced (evicted) in the D-cache
 - ♦ Prepares a modified block to be written back to lower memory
- Victim buffer decouples the write back to lower memory
 - ♦ Giving priority to read miss over write back to reduce miss penalty
- Problem: Victim buffer may hold block on a cache miss
 - ♦ Solution: transfer modified block in victim buffer into data cache



STUDENTS-HUB.com

Write Back with Write Allocate



STUDENTS-HUB.com

Write Through vs. Write Back



Unified vs. Separate Level 1 Cache

Unified Level 1 Cache (Princeton Memory Architecture). AKA Shared Cache

A single level 1 (L_1) cache is used for both instructions and data.

Separate instruction/data Level 1 caches (Harvard Memory Architecture):

The level 1 (L_1) cache is split into two caches, one for instructions (instruction cache, L_1 l-cache) and the other for data (data cache, L_1 D-cache).



Parallel vs. Serial Caches

Tag and Data usually separate (tag is smaller & faster)

- \diamond State bits stored along with tags
 - Valid bit, "LRU" bit(s), …

Parallel access to Tag and Data reduces latency (good for L1)



Serial access to Tag and Data reduces power (good for L2+)



Presentation Outline

Memory Hierarchy and the need for Cache Memory

The Basics of Caches

Multilevel Caches

Cache Performance and Memory Stall Cycles

Multilevel Caches

- Top level cache is kept small to
 - ♦ Reduce hit time
 - ♦ Reduce energy per access
- Add another cache level to
 - ♦ Reduce the memory gap
 - ♦ Reduce memory bus loading
- Multilevel caches can help
 - ♦ Reduce miss penalty
 - ♦ Reduce average memory access time
- Large L2 cache can capture many misses in L1 caches
- ♦ Reduce the global miss rate STUDENTS-HUB.com





Inclusion Policy

Multi-Level Cache Inclusion

- ♦ The content of the L1 cache is fully contained in the larger L2 cache
- ♦ Wastes L2 cache space, but L2 has space for additional blocks
- Multi-Level Cache Exclusion
 - ♦ The L2 cache can only have blocks that are not in the L1 cache
 - ♦ Prevents wasting space
- Non-Inclusive Non-Exclusion (NINE)
 - ♦ Multi-level inclusion and exclusion must be enforced by a protocol
 - \diamond If neither is enforced then L2 cache is neither inclusive nor exclusive

STUDENTS-HUB.com

Multilevel Inclusion

- L1 cache blocks are always present in L2 cache
 - ♦ Total number of unique cache blocks = L2 cache size
- ✤ A miss in L1, but a hit in L2 copies block from L2 to L1
- ✤ A miss in L1 and L2 brings a block into L1 and L2
- ✤ A write in L1 causes data to be written in L1 and L2
- Write-through policy is used from L1 to L2
- Write-back policy is used from L2 to lower-level memory
 - \diamond To reduce traffic on the memory bus
- ✤ A replacement (or invalidation) in L2 must be seen in L1

A block which is evicted from L2 must also be evicted from L1 STUDENTS-HUB.com
Uploaded By: Jibreel Bornat

Multilevel Exclusion

- L1 cache blocks are never found in L2 cache
 - ♦ Prevents wasting space
 - \diamond Total number of unique cache blocks = L1 + L2 cache size
- Cache miss in L1 and L2 brings the block into L1 only
- Cache miss in L1, but hit in L2 results in a swap of blocks
 - ♦ More complex to implement
- Block replaced in L1 is moved into L2
 - $\diamond\,$ L2 cache stores L1 evicted blocks, in case needed later in L1
 - ♦ L2 cache acts as a victim cache
- Write-Back policy from L1 to L2
- Write-Back policy from L2 to lower-level memory

STUDENTS-HUB.com

Non-Inclusive, Non Exclusive (NINE)

- L1 cache blocks may or may not be present in L2 cache
 - \diamond Total number of unique cache blocks > L2 cache size, but < (L1 + L2) size
- Cache miss in L1 and L2 brings a block into L1 and L2
- ✤ A miss in L1, but a hit in L2 copies block from L2 to L1
 - ♦ A block replaced in L1 can be discarded, if not modified
 - \diamond If a block is replaced and modified in L1, then it must be written back to L2
- ✤ A block can be replaced in L2 without removing it from L1
 - ♦ If L2 is inclusive, then block evicted from L2 must also be evicted from L1
- Write-back policy between L1 and L2

Write-through is not possible if a block exists in L1, but replaced in L2
 STUDENTS-HUB.com

Inclusive, Exclusive, and NINE L2 Cache



STUDENTS-HUB.com

Local and Global Miss Rates

Local Miss Rate

Number of cache misses / Memory accesses to this cache

Miss Rate_{L1} for L1 cache

Miss Rate_{L2} for L2 cache

Global Miss Rate

Number of cache misses / Memory accesses generated by processor

MissRate_{L1} for L1 cache (same as local)

MissRate_{L1} × **MissRate**_{L2} for L2 cache (different from local)

Global miss rate is a better measure for L2 cache

Fraction of the total memory accesses that miss in L1 and L2 STUDENTS-HUB.com Upload

Presentation Outline

- Memory Hierarchy and the need for Cache Memory
- The Basics of Caches
- Multilevel Caches
- Cache Performance and Memory Stall Cycles

Hit Rate and Miss Rate

- Hit Rate = Hits / (Hits + Misses)
- Miss Rate = Misses / (Hits + Misses)
- I-Cache Miss Rate = Miss rate in the Instruction Cache
- D-Cache Miss Rate = Miss rate in the Data Cache
- Example:
 - \diamond Out of 1000 instructions fetched, 60 missed in the I-Cache
 - \diamond 25% are load-store instructions, 50 missed in the D-Cache

 \diamond What are the I-cache and D-cache miss rates?

✤ I-Cache Miss Rate = 60 / 1000 = 6%

D-Cache Miss Rate = 50 / (25% × 1000) = 50 / 250 = 20%

STUDENTS-HUB.com

Memory Stall Cycles

- The processor stalls on a Cache miss
 - ♦ When fetching instruction from I-Cache and Block is not present
 - ♦ When loading/storing data in a D-cache and Block is not present
 - ♦ When writing data in a write-though D-cache and write buffer is full
- Miss Penalty: clock cycles to process a cache miss Miss Penalty is assumed equal for I-cache & D-cache Miss Penalty is assumed equal for Load and Store
- Memory Stall Cycles =
 - I-Cache Misses × Miss Penalty +
 - D-Cache Read Misses × Miss Penalty +
 - D-Cache Write Misses × Miss Penalty

STUDENTS-HUB.com

Combined Misses

- Combined Misses =
 - I-Cache Misses + D-Cache Read Misses + Write Misses
 - I-Cache Misses = I-Count × I-Cache Miss Rate
 - Read Misses = Load Count × D-Cache Read Miss Rate
 - Write Misses = Store Count × D-Cache Write Miss Rate
- Combined misses are often reported per 1000 instructions
- Memory Stall Cycles = Combined Misses × Miss Penalty

STUDENTS-HUB.com

Memory Stall Cycles Per Instruction

- Memory Stall Cycles Per Instruction =
 - Combined Misses Per Instruction × Miss Penalty
 - Miss Penalty is assumed equal for I-cache and D-cache
 - Miss Penalty is assumed equal for Load and Store
- Combined Misses Per Instruction =
 - I-Cache Miss Rate +
 - Load Frequency × D-Cache Read Miss Rate +
 - Store Frequency × D-Cache Write Miss Rate

STUDENTS-HUB.com

Example on Memory Stall Cycles

- Consider a program with the given characteristics
 - $\diamond~$ 20% of instructions are load and 10% are store
 - \diamond I-cache miss rate is 2%
 - $\diamond\,$ D-cache miss rate is 5% for load, and 1% for store
 - ♦ Miss penalty is 20 clock cycles for all cases (I-Cache & D-Cache)
 - ♦ Compute combined misses and stall cycles per instruction
- Combined misses per instruction in I-Cache and D-Cache
 - $2\% + 20\% \times 5\% + 10\% \times 1\% = 0.031$ combined misses per instruction
 - ♦ Equal to an average of 31 misses per 1000 instructions
- Memory stall cycles per instruction

 \Rightarrow 0.031 \times 20 (miss penalty) = 0.62 memory stall cycles per instruction

STUDENTS-HUB.com

CPU Time with Memory Stall Cycles

CPU Time =

(CPU execution cycles + Memory Stall Cycles) × Clock Cycle

CPU Time = I-Count × CPI_{overall} × Clock Cycle

CPI_{overall} = CPI_{execution} + Memory Stall Cycles per Instruction

CPI_{overall} = Overall CPI in the presence of cache misses

CPI_{execution} = Execution CPI (not counting cache misses)

Memory stall cycles per instruction increases the overall CPI

STUDENTS-HUB.com

Example on CPI with Memory Stalls

- ✤ A processor has CPI_{execution} = 1.5 (not counting cache misses)
 - \diamond I-Cache miss rate is 2%, D-cache miss rate is 5% for load & store
 - \diamond 20% of instructions are loads and stores
 - ♦ Cache miss penalty is 100 clock cycles for I-cache and D-cache
- What is the impact of cache misses on the overall CPI?
- Answer: Memory Stall Cycles per Instruction =

 0.02×100 (I-Cache) + $0.2 \times 0.05 \times 100$ (D-Cache) = 3

 $CPI_{overall} = 1.5 + 3 = 4.5$ cycles per instruction

 $CPI_{overall} / CPI_{execution} = 4.5 / 1.5 = 3$

Processor is 3 times slower due to memory stall cycles

STUDENTS-HUB.com

Average Memory Access Time

Average Memory Access Time (AMAT) AMAT = Hit time + Combined Miss rate × Miss penalty Hit Time = time to access the I-cache or D-cache (for a hit) Hit time is assumed to be the same for I-cache and D-cache

Combined Miss Rate for Instruction Access and Data Access

Combined Miss Rate = Combined Misses per Instruction Memory Accesses per Instruction

Memory Accesses per Instruction = 1 + %Load + %Store

Miss penalty is assumed to be the same for I-cache and D-cache STUDENTS-HUB.com Uploaded By: Jibreel Bornat

Memory Hierarchy Performance: Average Memory Access Time (AMAT), Memory Stall cycles

- The Average Memory Access Time (AMAT): The number of cycles required to complete an average memory access request by the CPU.
- Memory stall cycles per memory access: The number of stall cycles added to CPU execution cycles for one memory access.
- Memory stall cycles per average memory access = (AMAT -1)
- ✤ For ideal memory: AMAT = 1 cycle, this results in zero memory stall cycles.
- Memory stall cycles per average instruction =

Number of memory accesses per instruction

Instruction Fetch = (1 + fraction of loads/stores) x (AMAT -1)

Base $CPI = CPI_{execution} = CPI$ with ideal memory

CPI = CPI_{execution} + Mem Stall cycles per instruction

STUDENTS-HUB.com

AMAT Example

- Compute the overall average memory access time
 - ♦ Hit time = 1 clock cycle in both I-Cache and D-Cache
 - ♦ Miss penalty = 50 clock cycles for I-Cache and D-Cache
 - ♦ I-Cache misses = 3.8 misses per 1000 instructions
 - ♦ D-Cache misses = 41 misses per 1000 instructions
 - \diamond Load + Store frequency = 25%

Solution:

Combined Misses per Instruction = (3.8 + 41) / 1000 = 0.0448Combined Miss rate (per access) = 0.0448 / (1 + 0.25) = 0.03584Overall AMAT = $1 + 0.03584 \times 50 = 2.792$ cycles

STUDENTS-HUB.com

Cache Performance:(Ignoring Write Policy)Single Level L1 Princeton (Unified) Memory Architecture

CPUtime = Instruction count x CPI x Clock cycle time

CPI_{execution} = CPI with ideal memory

CPI = CPI_{execution} + Mem Stall cycles per instruction

Mem Stall cycles per instruction =

Memory accesses per instruction x Memory stall cycles per access

Assuming no stall cycles on a cache hit (cache access time = 1 cycle, stall = 0)Cache Hit Rate = H1Miss Rate = 1- H1i.e No hit penalty

Memory stall cycles per memory access = Miss rate x Miss penalty = (1-H1) x M AMAT = 1 + Miss rate x Miss penalty = 1 + (1-H1) x M Memory accesses per instruction = (1 + fraction of loads/stores) Miss Penalty = M = the number of stall cycles resulting from missing in cache = Main memory access time - 1

Thus for a unified L1 cache with no stalls on a cache hit:

CPI = CPI_{execution} + (1 + fraction of loads/stores) x (1 - H1) x M

AMAT = 1 + (1 - H1) x M STUDENTS-HUB.com

 CPI = CPI_{execution} + (1 + fraction of loads and stores) x stall cycles per access

 = CPI_{execution} + (1 + fraction of loads and stores) x (ANIAT + 1) breel Bornat

Memory Access Tree: For Unified Level 1 Cache

(Ignoring Write Policy)



AMAT = 1 + Stalls per average memory access

Cache Performance Example

- Suppose a CPU executes at Clock Rate = 200 MHz (5 ns per cycle) with a single level of cache.
- $CPI_{execution} = 1.1$
- Instruction mix: 50% arith/logic, 30% load/store, 20% control
- ✤ Assume a cache miss rate of 1.5% and a miss penalty of M= 50 cycles.

CPI = CPI_{execution} + mem stalls per instruction

Mem Stalls per instruction =

Mem accesses per instruction x Miss rate x Miss penalty Mem accesses per instruction = 1 + 0.3 = 1.3

Mem Stalls per memory access = $(1 - H1) \times M = 0.015 \times 50 = .75$ cycles

AMAT = 1 +.75 = 1.75 cycles

Mem Stalls per instruction = $1.3 \times 0.015 \times 50 = 0.975$

CPI = 1.1 + .975 = 2.075

STUDETINE ideal memory CPU with no misses is 2.075/1.1 = 1.88 times fastered Bornat

Cache Performance Example

- Suppose for the <u>previous example</u> we <u>double the clock rate</u> 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 = $M = 50 \times 2 = 100$ cycles.

 $CPI = 1.1 + 1.3 \times .015 \times 100 = 1.1 + 1.95 = 3.05$

Speedup =
$$(CPI_{old} \times C_{old})/(CPI_{new} \times C_{new})$$

= 2.075 x 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. STUDENTS-HUB.com

Cache Performance: Single Level L1 Harvard (Split) Memory Architecture



- For a CPU with separate or <u>split level one (L1)</u> caches for instructions and data (Harvard memory architecture) and no stalls for cache hits:
 - CPUtime = Instruction count x CPI x Clock cycle time
 - CPI = CPI_{execution} + Mem Stall cycles per instruction

Mem Stall cycles per instruction = Instruction Fetch Miss rate x M + Data Memory Accesses Per Instruction x Data Miss Rate x M

STUDENTS-HUB.com

Memory Access Tree For Separate Level 1 Caches



CPI = CPI_{execution} + Stall cycles per instruction = CPI_{execution} + (1 + fraction of loads/stores) x Stall Cycles per access

M = Miss Penalty = stall cycles per access resulting from missing in cache

M + 1 = Miss Time = Main memory access time

Data H1 = Level 1 Data Hit Rate1- Data H1 = Level 1 Data Miss Rate

Instruction H1 = Level 1 Instruction Hit Rate 1- Instruction H1 = Level 1 Instruction Miss Rate

% Instructions = Percentage or fraction of instruction fetches out of all memory accesses

% Data = Percentage or fraction of data accesses out of all memory accesses

Split L1 Cache Performance Example

- Suppose a CPU uses separate level one (L1) caches for instructions and data (Harvard memory architecture) with different miss rates for instruction and data access:
 - ♦ A cache hit incurs no stall cycles while a cache miss incurs 200 stall cycles for both memory reads and writes.
 - $\diamond \quad \mathbf{CPI}_{\text{execution}} = 1.1$
 - ♦ Instruction mix: 50% arith/logic, 30% load/store, 20% control
 - **Assume a cache miss rate of 0.5% for instruction fetch and a cache data miss rate of 6%.**
 - ♦ A cache hit incurs no stall cycles while a cache miss incurs 200 stall cycles for both memory reads and writes.
- Find the resulting stalls per access, AMAT and CPI using this cache?

CPI = CPI_{execution} + mem stalls per instruction

```
Memory Stall cycles per instruction = Instruction Fetch Miss rate x Miss Penalty +
Data Memory Accesses Per Instruction x Data Miss Rate x Miss Penalty
```

Memory Stall cycles per instruction = $0.5/100 \times 200 + 0.3 \times 6/100 \times 200 = 1 + 3.6 = 4.6$ cycles Stall cycles per average memory access = 4.6/1.3 = 3.54 cycles AMAT = 1 + 3.54 = 4.54 cycles

 $CPI = CPI_{execution} + mem stalls per instruction = 1.1 + 4.6 = 5.7 cycles$

What is the miss rate of a single level unified cache that has the same performance?

4.6 = 1.3 x Miss rate x 200 which gives a miss rate of 1.8 % for an equivalent unified cache

How much faster is the CPU with ideal memory?

The CPU with ideal cache (no misses) is 5.7/1.1 = 5.18 times faster

STUD Min from catal the CPI would have been = 1.1 + 1.3 X 200 = 261.1 cycles !! Uploaded By: Jibreel Bornat

Memory Access Tree For Separate Level 1 Caches Example

30% of all instructions executed are loads/stores, thus:

Fraction of instruction fetches out of all memory accesses = 1/(1+0.3) = 1/1.3 = 0.769 or 76.9 % Fraction of data accesses out of all memory accesses = 0.3/(1+0.3) = 0.3/1.3 = 0.231 or 23.1 %



STU Data = Percentage or fraction of data accesses out of all memory accesses = 23.1 %

Uploaded By: Jibreel Bornat

(Ignoring Write Policy)

Cache Write Miss Policy

Since data is usually not needed immediately on a write miss two options exist on a cache write miss:

(Bring old block to cache then update it)

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

No-Write Allocate:

Write Allocate:

i.e \underline{A} cache block frame is <u>not</u> allocated for the block to be modified (written-to)

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

While any of the above two write miss policies can be used with either write back or write through:

- <u>Write back caches always use write allocate</u> to capture subsequent writes to the block in cache.
- <u>Write through</u> caches <u>usually</u> use <u>no-write allocate</u> since subsequent writes still have to go to memory.

 Cache Write Miss
 = Block to be modified is not in cache

 STUDATIonate
 HABocate

 STUDATIONATE
 HABocate

 or assign a cache block frame for written data

Write Back Cache With Write Allocate: Cache Miss Operation



STUDENMS=+Miss_Penalty = stall cycles per access resulting from missing in cache

Memory Access Tree, Unified L₁ Write Through, No Write Allocate, <u>No Write Buffer</u>



Stall Cycles per access = AMAT - 1

M = Miss Penalty H1 = Level 1 Hit Rate 1- H1 = Level 1 Miss Rate

M = Miss Penalty = stall cycles per access resulting from missing in cacheM + 1 = Miss Time = Main memory access timeSTUDEN**STUDENI**+ U + U + Rate1 - H1 = Level 1 Miss Rate
Reducing Write Stalls For Write Though Cache Using Write Buffers

- To reduce write stalls when write though is used, a write buffer is used to eliminate or reduce write stalls:
 - Perfect write buffer: All writes are handled by write buffer, no stalling for writes

\diamond In this case (for unified L1 cache):

Stall Cycles Per Memory Access = % reads x (1 - H1) x M

(i.e No stalls at all for writes)

\diamond In this case (for unified L1 cache):

Stall Cycles/Memory Access = (% reads x (1 - H1) + % write stalls not eliminated) x M

STUDENTS-HUB.com

Write Through Cache Performance Example

- A CPU with $CPI_{execution} = 1.1$ Mem accesses per instruction = 1.3
- ✤ Uses a unified L1 <u>Write Through, No Write Allocate</u>, with:
 - ♦ No write buffer.

1

2

3

- ♦ Perfect Write buffer
- ♦ <u>A realistic write buffer</u> 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_{execution} + mem stalls per instruction % reads = 1.15/1.3 = 88.5% % writes = .15/1.3 = 11.5%
 - 1
 With No Write Buffer :
 Stall on all writes

 Mem Stalls/ instruction =
 1.3 x 50 x (88.5% x 1.5% + 11.5%) = 8.33 cycles

 CPI =
 1.1 + 8.33 =
 9.43
 - <u>With Perfect Write Buffer (all write stalls eliminated):</u>
 Mem Stalls/ instruction = 1.3 x 50 x (88.5% x 1.5%) = 0.86 cycles CPI = 1.1 + 0.86 = 1.96
- 3With Realistic Write Buffer (eliminates 85% of write stalls)
Mem Stalls/ instruction = $1.3 \times 50 \times (88.5\% \times 1.5\% + 15\% \times 11.5\%) = 1.98$ cycles
STUDENTS-HUB.comCPI = 1.1 + 1.98 = 3.08Uploaded By: Jibreel Bornat

Memory Access Tree Unified L₁ Write Back, With Write Allocate



AMAT = 1 + Stall Cycles Per Memory Access

CPI = **CPI**_{execution} + (1 + fraction of loads/stores) x Stall Cycles per access

M = Miss Penalty = stall cycles per access resulting from missing in cache

M + 1 = Miss Time = Main memory access time

H1 = Level 1 Hit Rate 1- H1 = Level 1 Miss Rate

STUDENTS-HUB.com

Write Back Cache Performance Example

- ✤ A CPU with CPI_{execution} = 1.1 uses a unified L1 with with <u>write back</u>, with <u>write</u> <u>allocate</u>, and the <u>probability a cache block is dirty = 10%</u>
- ✤ 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_{execution} + mem stalls per instruction Mem Stalls per instruction =

Mem accesses per instruction x Stalls per access Mem accesses per instruction = 1 + 0.3 = 1.3

Stalls per access = $(1-H1) \times (M \times \% \text{ clean} + 2M \times \% \text{ dirty})$

```
Stalls per access = 1.5\% \times (50 \times 90\% + 100 \times 10\%) = 0.825 cycles
AMAT = 1 + stalls per access = 1 + 0.825 = 1.825 cycles
Mem Stalls per instruction = 1.3 \times 0.825 = 1.07 cycles
CPI = 1.1 + 1.07 = 2.17
```

The ideal CPU with no misses is 2.17/1.1 = 1.97 times faster

STUDENTS-HUB.com

Memory Access Tree For Unified L₁

Write Back, With Write Allocate Example



M = Miss Penalty = 50 cyclesM + 1 = Miss Time = 50 + 1 = 51 cycles L1 access Time = 1 cycleSTUDEN**919350798.5%**1- H1 = 0.015 or 1.5%

Memory Access Tree Structure For Separate Level 1 Caches, Write Back, With Write Allocate



% Clean = Percentage or fraction of data L1 misses that are clean

ST

Dirty SPercentage or fraction of data L1 misses that are dirty = 1 - % Clean

2 Levels of Cache: L_1, L_2



L₁ = Level 1 Cache L₂ = Level 2 Cache Memory access penalty, M (stalls per main memory access) Access Time = M +1

Goal of multi-level Caches:

STI

Reduce the effective miss penalty incurred by level 1 cache misses by using additional levels of cache that capture some of these misses.

<u>Thus hiding more main memory latency and reducing AMAT further</u> DENTS-HUB com

Miss Rates For Multi-Level Caches

Local Miss Rate: This rate is the number of misses in a cache level divided by the number of memory accesses to this level (i.e those memory accesses that reach this level).

Local Hit Rate = 1 - Local Miss Rate

- Global Miss Rate: The number of misses in a cache level divided by the total number of memory accesses generated by the CPU.
- Since level 1 receives all CPU memory accesses, for level 1:

Local Miss Rate = Global Miss Rate = 1 - H1

For level 2 since it only receives those accesses missed in 1:

Local Miss Rate = Miss $rate_{L2} = 1 - H2$ Global Miss Rate = Miss $rate_{L1} \times Local Miss rate_{L2}$ = (1-H1) x (1 - H2)

STUDENTS-HUB.com

2-Level Cache (Both Unified) Performance (Ignoring Write Policy)

CPUtime = IC x (CPI_{execution} + Mem Stall cycles per instruction) x C

Mem Stall cycles per instruction = Mem accesses per instruction x Stall cycles per access

For a system with 2 levels of unified cache, assuming no penalty when found in L₁ cache: (T₁ = 0)

Stall cycles per memory access =

[miss rate L_1] x [Hit rate L_2 x Hit time L_2

+ Miss rate L₂ x Memory access penalty] =

$$(1-H1) \times H2 \times T2 + (1-H1)(1-H2) \times M$$

$$L1 \text{ Miss, L2 Hit}$$

$$H1 = L1 \text{ Hit Rate}$$

$$T1 = \text{stall cycles per L1 access hit}$$

$$H2 = \text{Local L2 Hit Rate}$$

$$T2 = \text{stall cycles per L2 access hit}$$

$$Full \text{ Miss}$$

$$Full \text{ Miss}$$

= $CPI_{execution}$ + (1 + fraction of loads and stores) x (AMAT - 1)

2-Level Cache (Both Unified) Performance Memory Access Tree (Ignoring Write Policy) CPU Stall Cycles Per Memory Access



STUDENTS-HUB.com

Unified Two-Level Cache Example

- CPU with $CPI_{execution} = 1.1$ running at clock rate = 500 MHz
- ✤ 1.3 memory accesses per instruction.
- With two levels of cache (both unified)
- ↔ L_1 hit access time = 1 cycle (no stall on a hit, T1= 0), a miss rate of 5%
- ↔ L_2 hit access time = 3 cycles (T2= 2 stall cycles per hit) with local miss rate 40%,
- ✤ Memory access penalty, M = 100 cycles (stalls per access). Find CPI ...

 $\begin{array}{rcl} \mathsf{CPI} = & \mathsf{CPI}_{\mathsf{execution}} + \mathsf{Mem Stall cycles per instruction} \\ \text{With No Cache, } & \mathsf{CPI} = 1.1 + 1.3 \times 100 = 131.1 \\ \text{With single L}_1, & \mathsf{CPI} = 1.1 + 1.3 \times .05 \times 100 = 7.6 \\ \text{Mem Stall cycles per instruction} = & \mathsf{Mem accesses per instruction} \times \mathsf{Stall cycles per access} \\ \text{Stall cycles per memory access} = & (1-H1) \times H2 \times T2 & + & (1-H1)(1-H2) \times M \\ & = & 0.05 \times .6 \times 2 & + & 0.05 \times 0.4 \times 100 \\ & = & 0.06 + & 2 & = & 2.06 \text{ cycles} \end{array}$

AMAT = 2.06 + 1 = 3.06 cycles

Mem Stall cycles per instruction = Mem accesses per instruction x Stall cycles per access

= $2.06 \times 1.3 = 2.678$ cycles CPI = 1.1 + 2.678 = 3.778Speedup = 7.6/3.778 = 2

STUDENTS-HUB.com

CPI = CPI_{execution} + (1 + fraction of loads and stores) x stall cycles per access = CPI_{execution} + (1 + fraction of loads and stores Bx/(AMANee1)Bornat

(Ignoring Write Policy)

Memory Access Tree For 2-Level Cache (Both Unified) Example CPU Stall Cycles Per Memory Access



STUDE NTS = $COPB_{execution}$ + (1 + fraction of loads and stores) x (AMAT - 1)

Memory Access Tree Structure For 2-Level Cache (Separate Level 1 Caches, Unified Level 2)



Common Write Policy For 2-Level Cache

- L₁ **Write Policy For Level 1 Cache:**
 - ♦ Usually <u>Write through to Level 2</u>.

(not write through to main memory just to L2)

- ♦ Write allocate is used to reduce level 1 read misses.
- \diamond Use <u>write buffer</u> to reduce write stalls to level 2.
- L₂ ♦ Write Policy For Level 2 Cache:
 - ♦ Usually write back with write allocate is used.
 - To minimize memory bandwidth usage.
 - The above 2-level cache write policy results in <u>inclusive L2 cache</u> since the content of L1 is also in L2
 - Common in the majority of all CPUs with 2-levels of cache
 - As opposed to exclusive L1, L2 (e.g AMD Athlon XP, A64)



STUDENTS-HUB.com

2-Level (Both Unified) Memory Access Tree L1: Write Through to L2, Write Allocate, With Perfect Write Buffer L2: Write Back with Write Allocate



AMAT = 1 + Stall Cycles Per Memory AccessSTUCPE A CPI HUB. #Q(I) + fraction of loads and stores) x Stall Cycles per access

Two-Level (Both Unified) Cache Example With Write Policy

- ✤ CPU with CPI_{execution} = 1.1 running at clock rate = 500 MHz
- ✤ 1.3 memory accesses per instruction. Two levels of cache (both unified)
- ✤ For L₁:
 - ♦ Čache operates at 500 MHz (no stall on L1 Hit, T1 =0) with a miss rate of 1-H1 = 5%
 - \diamond Write though to L₂ with perfect write buffer with write allocate
- ✤ For L₂:
 - \Rightarrow Hit access time = 3 cycles (T2= 2 stall cycles per hit) local miss rate 1- H2 = 40%
 - ♦ Write back to main memory with write allocate
 - ♦ Probability a cache block is dirty = 10%
- ✤ Memory access penalty, M = 100 cycles.
- Create memory access tree and find, stalls per memory access, AMAT, CPI.
- ✤ Stall cycles per memory access = (1-H1) x H2 x T2 +

 $(1 - H1) \times (1 - H2) \times (\% \text{ clean } \times M + \% \text{ dirty } \times 2M)$ = .05 x .6 x 2 + .05 x .4 x (.9 x 100 + .1 x200) = .06 + 0.02 x 110 = .06 + 2.2 = 2.26

✤ AMAT = 2.26 + 1 = 3.26 cycles

Mem Stall cycles per instruction = Mem accesses per instruction x Stall cycles per access

CPI = 1.1 + 2.938 = 4.038 = 4

 $STUDERESCPI _ SCPI _ CPI _ C$

Memory Access Tree For Two-Level (Both Unified) Cache Example With Write Policy

L1: Write Through to L2, Write Allocate, With Perfect Write Buffer L2: Write Back with Write Allocate



AMAT = 1 + Stall Cycles Per Memory AccessSTUCPE ACE HUB. #0(I) + fraction of loads and stores) x Stall Cycles per access

Memory Access Tree Structure For 2-Level Cache(Separate Level 1 Caches, Unified Level 2) L1: Write Through to L2, Write Allocate, With Perfect Write Buffer L2: Write Back with Write Allocate

