



中国科学技术大学
University of Science and Technology of China

计算机体系结构

周学海

xhzhou@ustc.edu.cn

0551-63606864

中国科学技术大学



Review: Cache基本概念

- **存储系统分层结构**
 - 为什么要这样做?
 - 为什么可以这样做?
- **存储系统的性能指标：速度、容量和价格**
 - 平均访存时间=命中时间+失效率×失效开销
- **Cache 4Q**
 - 映射规则
 - 查找方法
 - 替换策略
 - 写策略



第4章 存储层次结构设计

4.1 Cache的基本概念

存储系统的层次结构

Cache基本知识

4.2 Cache的基本优化方法

4.3 Cache的高级优化方法

4.4 存储器技术与优化

4.5 虚拟存储器 - 基本原理



4.2 Cache 基本优化方法

Cache
性能分析

降低
失效率

减少
失效开销

缩短
命中时间

平均访存时间 = 命中时间 + 失效率 × 失效开销



Cache 性能分析

- **CPU time = (CPU execution clock cycles + Memory stall clock cycles) x clock cycle time**
- **Memory stall clock cycles = (Reads x Read miss rate x Read miss penalty + Writes x Write miss rate x Write miss penalty)**
- **Memory stall clock cycles = Memory accesses x Miss rate x Miss penalty**
- **Different measure: AMAT**

Average Memory Access time (AMAT) = Hit Time + (Miss Rate x Miss Penalty)

- **Note: *memory hit time is included in execution cycles.***



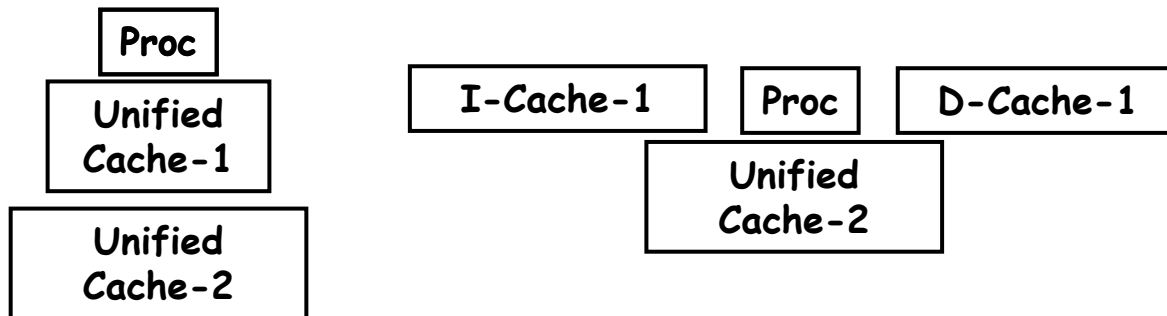
性能分析举例

- **Suppose a processor executes at**
 - Clock Rate = 200 MHz (5 ns per cycle), Ideal (no misses) CPI = 1.1
 - 50% arith/logic, 30% ld/st, 20% control
- **Miss Behavior:**
 - 10% of memory operations get 50 cycle miss penalty
 - 1% of instructions get same miss penalty
- **CPI = ideal CPI + average stalls per instruction**
$$= 1.1(\text{cycles/ins}) + [0.30 (\text{DataMops/ins}) \times 0.10 (\text{miss/DataMop}) \times 50 (\text{cycle/miss})] + [1 (\text{InstMop/ins}) \times 0.01 (\text{miss/InstMop}) \times 50 (\text{cycle/miss})]$$
$$= (1.1 + 1.5 + .5) \text{ cycle/ins} = 3.1$$
- **65% (2/3.1) of the time the proc is stalled waiting for memory!**
- **AMAT = $(1/1.3) \times [1 + 0.01 \times 50] + (0.3/1.3) \times [1 + 0.1 \times 50] = 2.54$**



Example: Harvard Architecture

- **Unified vs Separate I&D (Harvard)**



- **Statistics (given in H&P):**

- 16KB I&D: Inst miss rate=0.64%, Data miss rate=6.47%
- 32KB unified: Aggregate miss rate=1.99%

- **Which is better (ignore L2 cache)?**

- Assume 33% data ops \Rightarrow 75% accesses from instructions (1.0/1.33)
- hit time=1, miss time=50
- Note that data hit has 1 stall for unified cache (only one port)

- **$AMATHarvard = 75\% \times (1 + 0.64\% \times 50) + 25\% \times (1 + 6.47\% \times 50) = 2.05$**

- **$AMATUnified = 75\% \times (1 + 1.99\% \times 50) + 25\% \times (1 + 1 + 1.99\% \times 50) = 2.24$**



Size	Instruction cache	Data cache	Unified cache
8 KB	8.16	44.0	63.0
16 KB	3.82	40.9	51.0
32 KB	1.36	38.4	43.3
64 KB	0.61	36.9	39.4
128 KB	0.30	35.3	36.2
256 KB	0.02	32.6	32.9

FIGURE 5.8 Miss per 1000 instructions for instruction, data, and unified caches of different sizes. The percentage of instruction references is about 78%. The data are for two-way associative caches with 64-byte blocks for the same computer and benchmarks as Figure 5.6.



- 以顺序执行的计算机 UltraSPARC III为例. 假设Cache失效开销为100 clock cycles, 所有指令忽略存储器停顿需要1个cycle, Cache失效可以用两种方式给出

(1) 假设平均失效率为2%, 平均每条指令访存1.5次

(2) 假设每1000条指令cache失效次数为30次

分别基于上述两种条件计算处理器的性能

- 结论:

- $$\text{CPUtime} = \text{IC} * (1 + 2\% * 1.5 * 100) * T = \text{IC} * 4 * T$$

(1) CPI越低, 固定周期数的Cache失效开销的相对影响就越大

(2) 在计算CPI时, 失效开销的单位是时钟周期数。因此, 即使两台计算机的存储层次完全相同, 时钟频率较高的CPU的失效开销会较大, 其CPI中存储器停顿部分也就较大。

因此 Cache对于低CPI, 高时钟频率的CPU来说更加重要



考虑不同组织结构的Cache对性能的影响:

- **B-19例题: 直接映像Cache 和两路组相联Cache, 试问他们对CPU性能的影响? 先求平均访存时间, 然后再计算CPU性能。分析时请用以下假设:**
 - (1) **理想Cache(命中率为100%) 情况下CPI 为1.0, 时钟周期为0.35ns, 平均每条指令访存1.4次**
 - (2) **两种Cache容量均为128KB, 块大小都是64B**
 - (3) **采用组相联时, 由于多路选择器的存在, 时钟周期增加到原来的1.35倍**
 - (4) **两种结构的失效开销都是65ns (在实际应用中, 应取整为整数个时钟周期)**
 - (5) **命中时间为1个cycle, 128KB直接映像Cache的失效率为2.1%, 相同容量的两路组相联Cache 的失效率为1.9%**



失效开销与Out-of-Order执行的处理器

$$\frac{MemoryStallCycles}{Instruction} = \frac{Misses}{Instruction} \times (TotalMissLatency - OverlappedMissLatency)$$

- 需要确定两个参数：
- Length of memory latency
- Length of latency overlap
- 例如：在前面的例子中，若假设允许处理器乱序执行，则对于直接映射方式，假设30%的失效开销可以覆盖 (overlap)，那么原来的70ns失效开销就变为49ns.



Summary of performance equations in this chapter

$$2^{\text{index}} = \frac{\text{Cache size}}{\text{Block size} \times \text{Set associativity}}$$

$$\text{CPU execution time} = (\text{CPU clock cycles} + \text{Memory stall cycles}) \times \text{Clock cycle time}$$

$$\text{Memory stall cycles} = \text{Number of misses} \times \text{Miss penalty}$$

$$\text{Memory stall cycles} = \text{IC} \times \frac{\text{Misses}}{\text{Instruction}} \times \text{Miss penalty}$$

$$\frac{\text{Misses}}{\text{Instruction}} = \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}}$$

$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

$$\text{CPU execution time} = \text{IC} \times \left(\text{CPI}_{\text{execution}} + \frac{\text{Memory stall clock cycles}}{\text{Instruction}} \right) \times \text{Clock cycle time}$$

$$\text{CPU execution time} = \text{IC} \times \left(\text{CPI}_{\text{execution}} + \frac{\text{Misses}}{\text{Instruction}} \times \text{Miss penalty} \right) \times \text{Clock cycle time}$$

$$\text{CPU execution time} = \text{IC} \times \left(\text{CPI}_{\text{execution}} + \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}} \times \text{Miss penalty} \right) \times \text{Clock cycle time}$$

$$\frac{\text{Memory stall cycles}}{\text{Instruction}} = \frac{\text{Misses}}{\text{Instruction}} \times (\text{Total miss latency} - \text{Overlapped miss latency})$$

$$\text{Average memory access time} = \text{Hit time}_{L1} + \text{Miss rate}_{L1} \times (\text{Hit time}_{L2} + \text{Miss rate}_{L2} \times \text{Miss penalty}_{L2})$$

$$\frac{\text{Memory stall cycles}}{\text{Instruction}} = \frac{\text{Misses}_{L1}}{\text{Instruction}} \times \text{Hit time}_{L2} + \frac{\text{Misses}_{L2}}{\text{Instruction}} \times \text{Miss penalty}_{L2}$$

Figure B.7 Summary of performance equations in this appendix. The first equation calculates the cache index size, and the rest help evaluate performance. The final two equations deal with multilevel caches, which are explained early in the next section. They are included here to help make the figure a useful reference.



改进Cache 性能的方法

- **平均访存时间 = 命中时间 + 失效率 × 失效开销**
- **从上式可知，基本途径**
 - 降低失效率
 - 减少失效开销
 - 缩短命中时间



Cache基本优化方法

- **降低失效率**
 - 1、增加Cache块的大小
 - 2、增大Cache容量
 - 3、提高相联度
- **减少失效开销**
 - 4、多级Cache
 - 5、使读失效优先于写失效
- **缩短命中时间**
 - 6、避免在索引缓存期间进行地址转换



4.2 Cache 基本优化方法

Cache
性能分析

降低
失效率

减少
失效开销

缩短
命中时间

$$\text{平均访存时间} = \text{命中时间} + \text{失效率} \times \text{失效开销}$$



降低失效率

Cache失效的原因 可分为三类 3C

- **强制性失效 (Compulsory)**

- 第一次访问某一块，只能从下一级Load，也称为冷启动或首次访问失效

- **容量失效 (Capacity)**

- 如果程序执行时，所需块由于容量不足，不能全部调入Cache，则当某些块被替换后，若又重新被访问，就会发生失效。
- 可能会发生“抖动”现象

- **冲突失效 (Conflict (collision))**

- 组相联和直接相联的副作用
- 若太多的块映像到同一组（块）中，则会出现该组中某个块被别的块替换（即使别的组或块有空闲位置），然后又被重新访问的情况，这就属于冲突失效



各种类型的失效率

Compulsory misses are independent of cache size

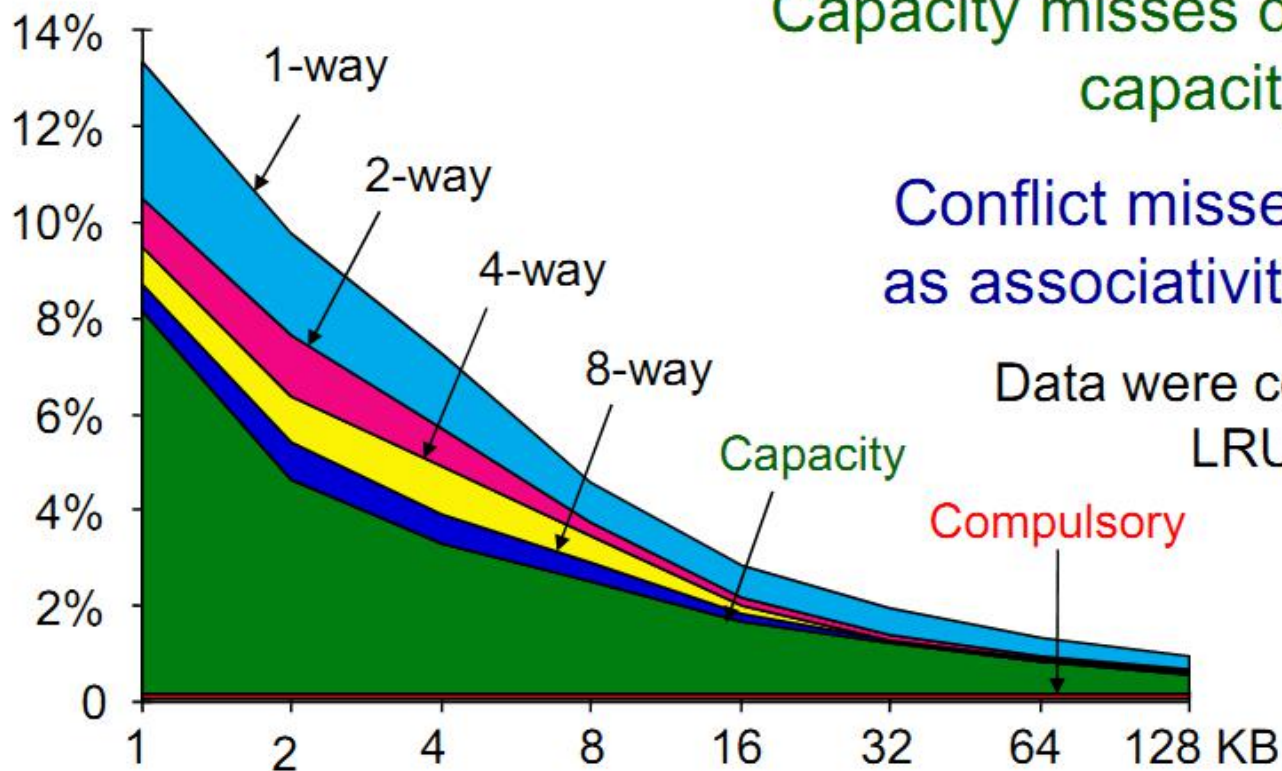
Very small for long-running programs

Capacity misses decrease as capacity increases

Conflict misses decrease as associativity increases

Data were collected using LRU replacement

Miss Rate



Compulsory



从统计规律中得到的一些结果

- 相联度越高，冲突失效就越小
- 强制性失效和容量失效不受相联度的影响
- 强制性失效不受Cache容量的影响
- 容量失效随着容量的增加而减少
- **符合2:1 Cache经验规则**
 - 即大小为N的直接映象Cache的失效率约等于大小为N/2的两路组相联的Cache失效率。



减少3C的方法

从统计规律可知

- **增大Cache容量**

- 对冲突和容量失效的减少有利

- **增大块**

- 减缓强制性失效

- 可能会增加冲突失效（因为在容量不变的情况下，块的数目减少了）

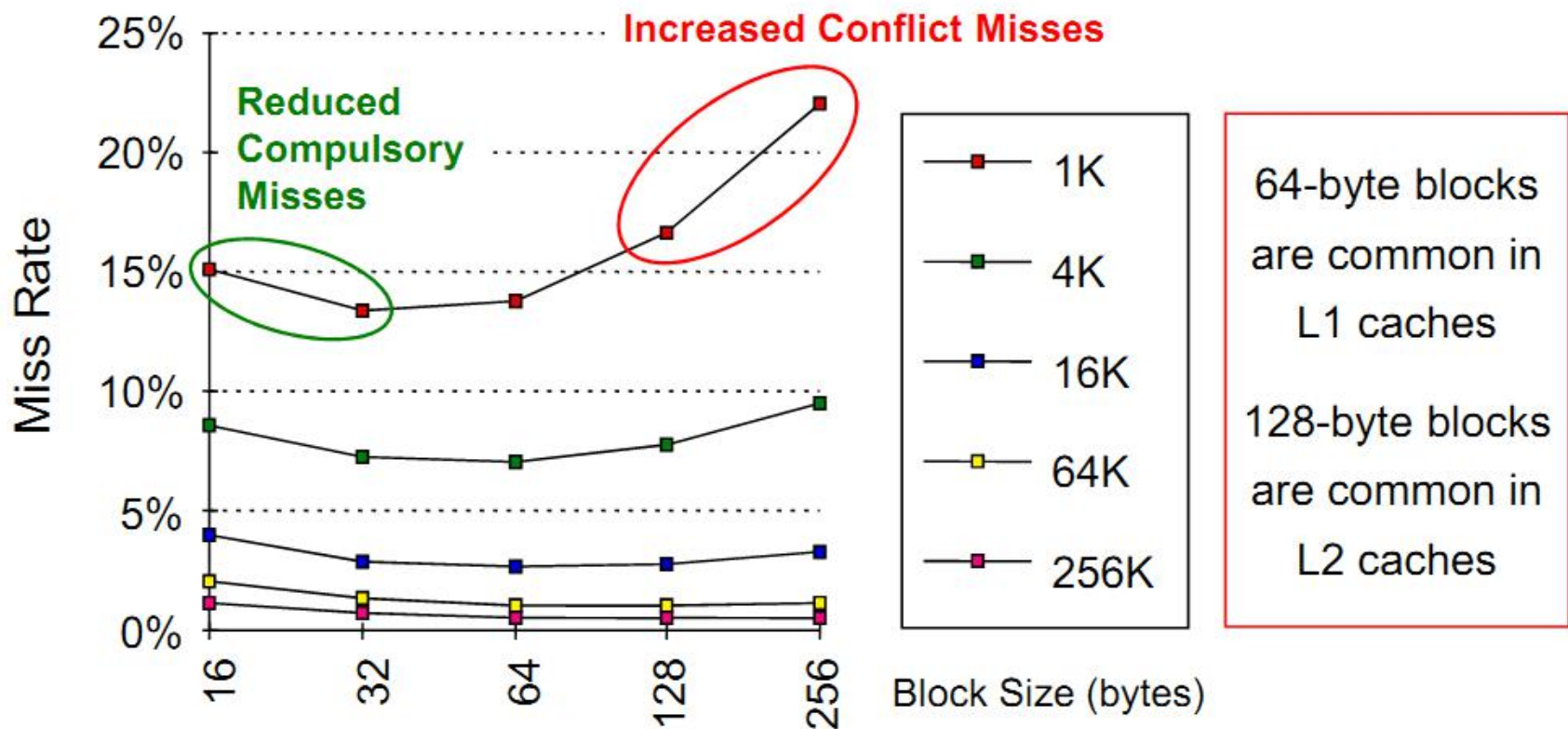
- **通过预取可帮助减少强制性失效**

- 必须小心不要把你需要的东西换出去

- 需要预测比较准确（对数据较困难，对指令相对容易）



增加块的大小





块大小、容量对失效率的影响

- 已知: 不同Cache容量以及块大小条件下的失效率

Block Size	Cache = 4 KB	Cache = 16 KB	Cache = 64 KB	Cache = 256 KB
16 bytes	8.57%	3.94%	2.04%	1.09%
32 bytes	7.24%	2.87%	1.35%	0.70%
64 bytes	7.00%	2.64%	1.06%	0.51%
128 bytes	7.78%	2.77%	1.02%	0.49%
256 bytes	9.51%	3.92%	1.15%	0.49%

- **Memory latency = 80 cycles + 1 cycle per 8 bytes**
 - Latency of 16-byte block = $80 + 2 = 82$ clock cycles
 - Latency of 32-byte block = $80 + 4 = 84$ clock cycles
 - Latency of 256-byte block = $80 + 32 = 112$ clock cycles
- **Which block has smallest AMAT for each cache size?**



块大小、容量对AMAT的影响

- **Solution: assume hit time = 1 clock cycle**
 - Regardless of block size and cache size
- **Cache Size = 4 KB, Block Size = 16 bytes**
 - $AMAT = 1 + 8.57\% \times 82 = 8.027$ clock cycles
- **Cache Size = 256 KB, Block Size = 256 bytes**
 - $AMAT = 1 + 0.49\% \times 112 = 1.549$ clock cycles

Block Size	Cache = 4 KB	Cache = 16 KB	Cache = 64 KB	Cache = 256 KB
16 bytes	AMAT = 8.027	AMAT = 4.231	AMAT = 2.673	AMAT = 1.894
32 bytes	AMAT = 7.082	AMAT = 3.411	AMAT = 2.134	AMAT = 1.588
64 bytes	AMAT = 7.160	AMAT = 3.323	AMAT = 1.933	AMAT = 1.449
128 bytes	AMAT = 8.469	AMAT = 3.659	AMAT = 1.979	AMAT = 1.470
256 bytes	AMAT = 11.65	AMAT = 4.685	AMAT = 2.288	AMAT = 1.549



块大小、容量对AMAT的影响

Block Size	Penalty	Cache Size 1K	Cache Size 4K	Cache Size 16K	Cache Size 64K	Cache Size 256K
16	42	15.05% / 7.321	8.57% / 4.599	3.94% / 2.655	2.04% / 1.857	1.09% / 1.458
32	44	13.34% / 6.870	7.24% / 4.186	2.87% / 2.263	1.35% / 1.594	0.70% / 1.308
64	48	11.76% / 7.605	7.00% / 4.360	2.64% / 2.267	1.06% / 1.509	0.51% / 1.245
128	56	16.64% / 10.318	7.78% / 5.357	2.77% / 2.551	1.02% / 1.571	0.49% / 1.274
256	72	22.01% / 16.847	9.51% / 7.847	3.29% / 3.369	1.15% / 1.828	0.49% / 1.353

Miss Rate / Average Access Time (in cycles)

Remember

what's going on here?

• $\text{Avg-access-time} = \text{hit-time} + \text{miss-rate} \times \text{miss-penalty}$

- 降低失效率最简单的方法是增加块大小；统计结果如图所示
- 假定存储系统在延迟40个时钟周期后，每2个时钟周期能送出16个字节，即：经过42个时钟周期，它可提供16个字节；经过44个四周周期，可提供32个字节；依此类推。试根据图5-6列出的各种容量的Cache，在块大小分别为多少时，平均访存时间最小？



块大小、容量的权衡

- **从统计数据可得到如下结论**

- 对于给定Cache容量，块大小增加时，失效率开始是下降，但后来反而上升
- Cache容量越大，使失效率达到最低的块大小就越大

- **分析**

- 块大小增加，可使强制性失效减少（空间局部性原理）
- 块大小增加，可使冲突失效增加（Cache中块数量减少）
- 失效开销增大（上下层间移动，数据传输时间变大）

- **设计块大小的原则，不能仅看失效率**

- 原因：平均访存时间 = 命中时间 + 失效率 × 失效开销



提高相联度

- 8路组相联在降低失效率方面的作用已经和全相联一样有效
- **2:1 Cache经验规则**
 - 容量为N的直接映象Cache失效率与容量为N/2的两路组相联Cache的失效率差不多相同
- **提高相联度，会增加命中时间**

Size (KB)	1-way	2-way	4-way	8-way
1	7.65	6.60	6.22	5.44
2	5.90	4.90	4.62	4.09
4	4.60	3.95	3.57	3.19
8	3.30	3.00	2.87	2.59
16	2.45	2.20	2.12	2.04
32	2.00	1.80	1.77	1.79
64	1.70	1.60	1.57	1.59
128	1.50	1.45	1.42	1.44

Average Memory Access Time

OOPS!

Victim Cache(1/2)

- 在Cache和Memory之间增加一个小的全相联Cache

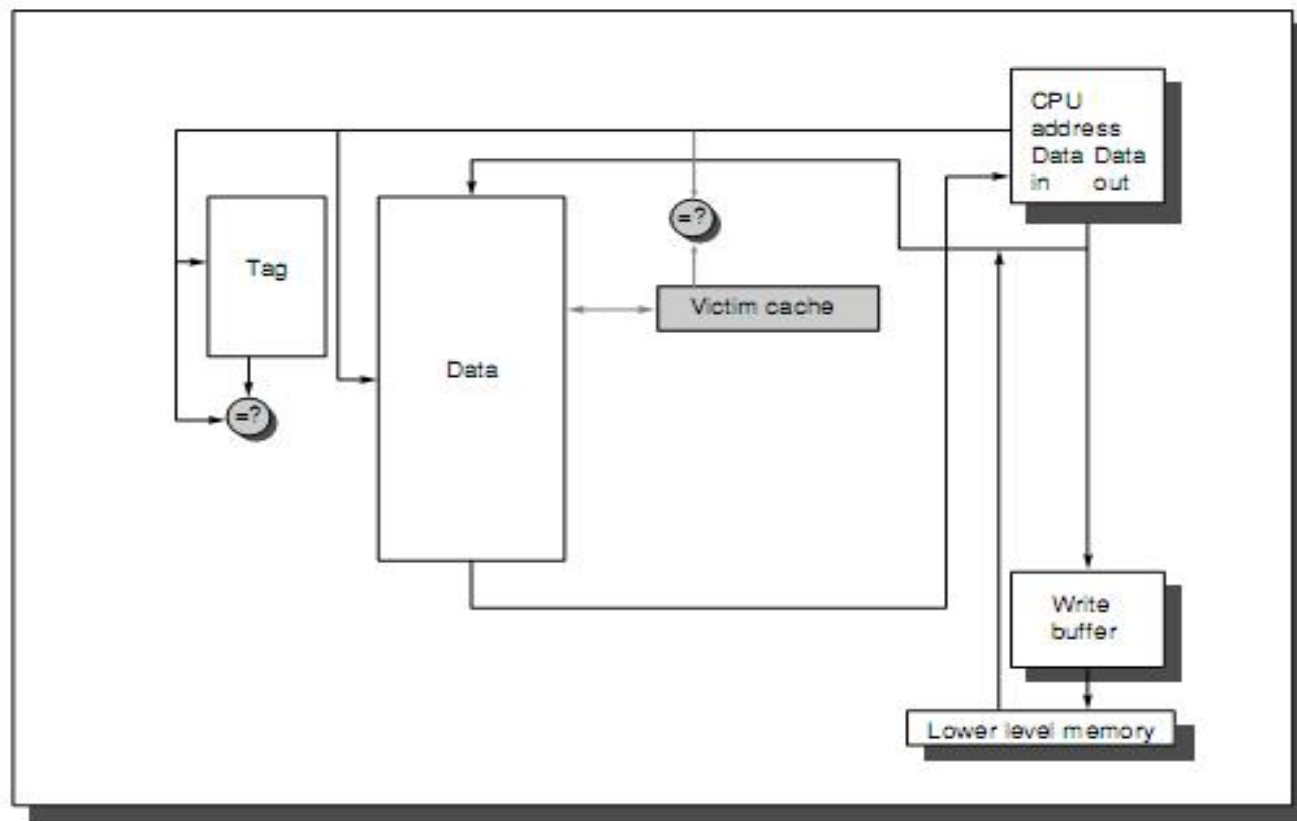


FIGURE 5.13 Placement of victim cache in the memory hierarchy. Although it reduces miss penalty, the victim cache is aimed at reducing the damage done by conflict misses, described in the next section. Jouppi [1990] found the four-entry victim cache could reduce the miss penalty for 20% to 95% of conflict misses.



Victim Cache(2/2)

- **基本思想**

- 通常Cache为直接映象时冲突失效率较大
- Victim cache采用全相联 - 失效率较低
- Victim cache存放由于（冲突）失效而被丢弃的那些块
- 失效时，首先检查Victim cache是否有该块，如果有就将该块与Cache中相应块比较。

- **Jouppi (DEC SRC)发现**

- 含1到5项的Victim cache对减少失效很有效，尤其是对于那些小型的直接映象数据Cache。
- 测试结果，项为4的Victim Cache能使4KB直接映象数据Cache冲突失效减少20%-90%



4.2 Cache 基本优化方法

Cache
性能分析

降低
失效率

减少
失效开销

缩短
命中时间

平均访存时间 = 命中时间 + 失效率 × 失效开销



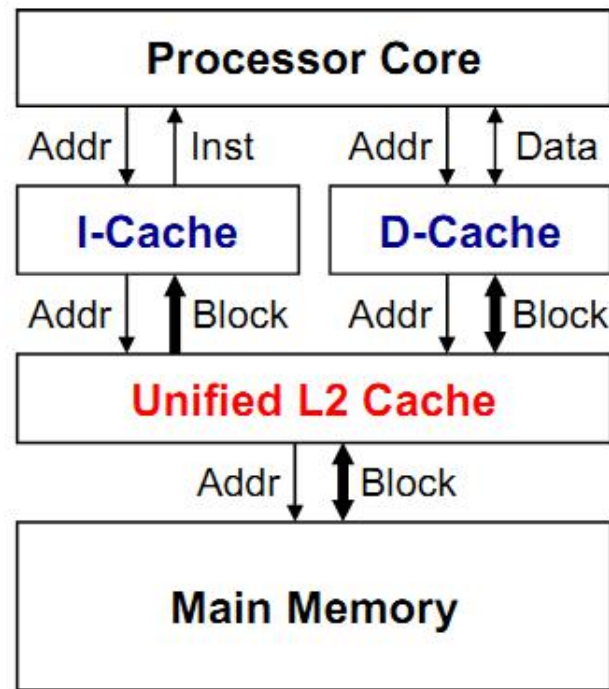
减少失效开销

- **减少CPU与存储器间性能差异的重要手段**
 - 平均访存时间 = 命中时间 + 失效率 × 失效开销
- **基本手段:**
 - 4、多级Cache技术(Multilevel Caches)
 - 5、让读优先于写(Giving Priority to Read Misses over Writes)



采用多级Cache

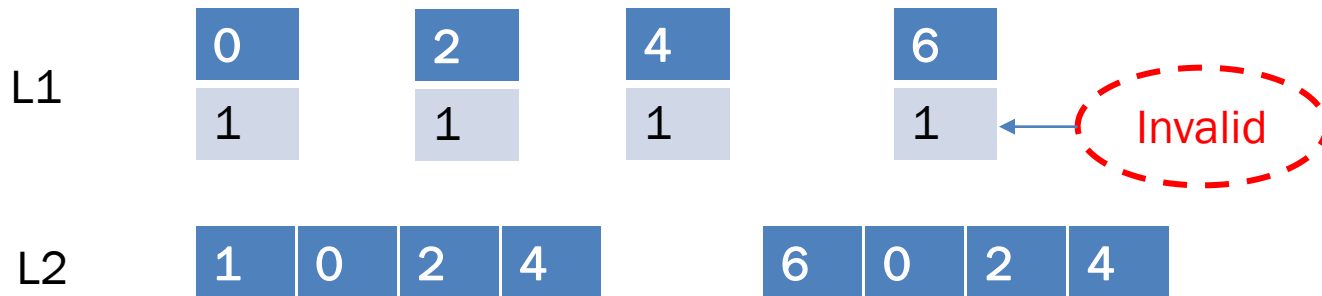
- **一级cache保持较小容量**
 - 降低命中时间
 - 降低每次访问的能耗
- **增加二级cache**
 - 减少与存储器的gap
 - 减少存储器总线的负载
- **多级cache的优点**
 - 减少失效开销
 - 缩短平均访存时间 (AMAT)
- **较大容量的L2 cache可以捕捉许多L1 cache的失效**
 - 降低全局失效率





多级包容性(multilevel inclusive)

- **L1 cache 的块总是存在于L2 cache中**
 - 浪费了L2 cache 空间, L2 还应当有存放其他块的空间
- **L1中miss, 但在L2中命中, 则从L2拷贝相应的块到L1**
- **在L1和L2中均miss, 则从更低级拷贝相应的块到L1和L2**
- **对L1写操作导致将数据同时写到L1和L2**
- **Write-through 策略用于L1到L2**
- **Write-back 策略可用于L2 到更低级存储器, 以降低存储总线的数据传输压力**
- **L2的替换动作 (或无效)对L1可见**
 - 即L2的一块被替换出去, 那么其在L1中对应的块也要被替换出去。
- **L1和L2的块大小可以相同也可以不同**
 - Pentium 4 had 64-byte blocks in L1 but 128-byte blocks in L2
 - Core i7 uses 64-byte blocks at all cache levels (simpler)



- L1: 1-way, 容量: 2块
- L2 4-way, 容量: 4块
- 块大小相同
- 访问的块序列10**2**4**6**
- 替换策略: LRU



0
1

4
1

8
1

Invalid

0(0)	
0(1)	

0(0)	
0(1)	
2(4)	
2(5)	

4(8)	
4(9)	
2(4)	
2(5)	

0	0(0)
1	0(1)
2	1(2)
3	1(3)
4	2(4)
5	2(5)
6	3(6)
7	3(7)
8	4(8)
9	4(9)

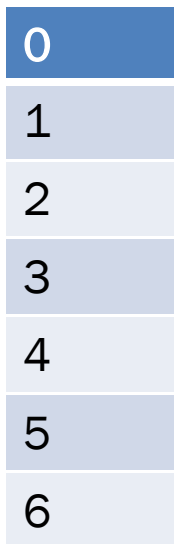
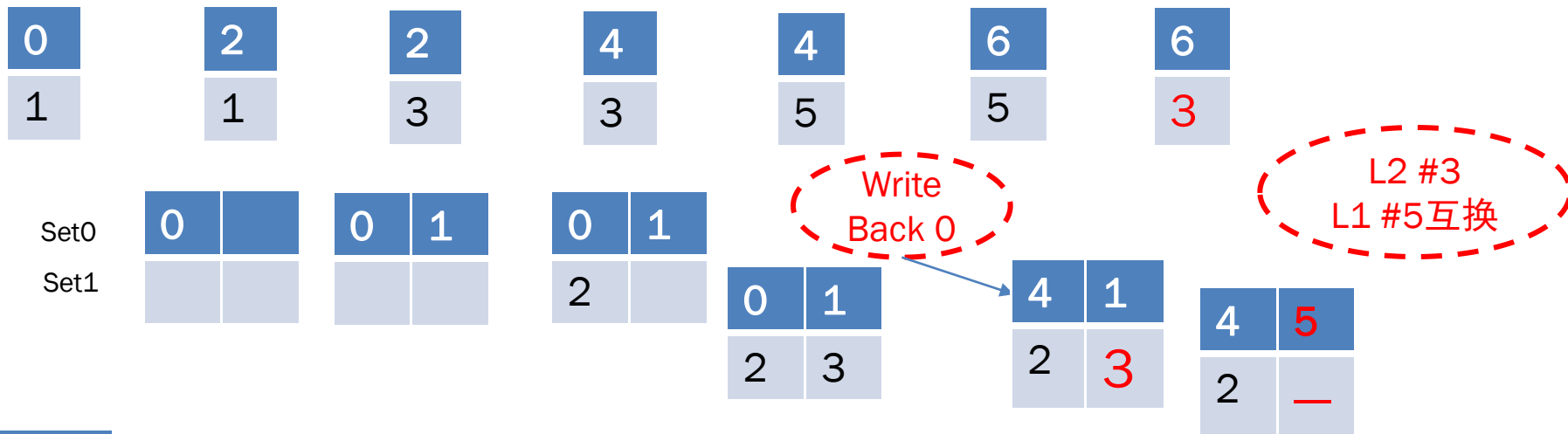
- L1: 1-way, 容量: 2块
- L2: 2-way, 容量: 4块
- L1和L2的块大小不同, L1的blockSize是L2的blockSize的1/2
- 替换策略: LRU
- 考察按L1的块划分的访问序列:

0148



多级不包容 (Multilevel Exclusive)

- **L1 cache 中的块不会在L2 cache中，以避免浪费空间**
- **在L1中miss, 但在L2中命中，将导致Cache间块的互换**
- **在L1和L2均miss, 将仅仅从更低层拷贝相应的块到L1**
- **L1的被替换的块移至L2**
 - L2 存储L1抛弃的块，以防后续L1还需要使用
- **L1到L2的写策略为 Write-Back**
- **L2到更低级cache的写策略为 Write-Back**
- **AMD Athlon: 64KB L1、256KB L2**



- L1: 1-way L2: 2-way
- 访问的块序列 0 1 2 3 4 5 6 3
- 替换策略: LRU
- 块大小相同



多级cache的性能分析

- **局部失效率**: 该级Cache的失效次数 / 到达该级Cache的访存次数
 - Miss rateL1 for L1 cache
 - Miss rateL2 for L2 cache
- **全局失效率**: 该级Cache的失效次数/ CPU发出的访存总次数
 - Miss rateL1 for L1 cache
 - Miss rateL1 × Miss rateL2 for L2 cache
 - 全局失效率是度量L2 cache性能的更好方法
- **性能参数**
 - $AMAT = Hit\ Time_{L1} + Miss\ rate_{L1} \times Miss\ penalty_{L1}$
 - $Miss\ penalty_{L1} = Hit\ Time_{L2} + Miss\ rate_{L2} \times Miss\ penalty_{L2}$
 - $AMAT = Hit\ Time_{L1} + Miss\ rate_{L1} \times (Hit\ Time_{L2} + Miss\ rate_{L2} \times Miss\ penalty_{L2})$



L1 cache 失效率

- **对于I-Cache和D-Cache分开的L1 Cache**

$$\text{Miss Rate}_{L1} = \%inst \times \text{Miss Rate}_{I\text{-Cache}} + \%data \times \text{Miss Rate}_{D\text{-Cache}}$$

$$\%inst = \text{Percent of Instruction Accesses} = 1 / (1 + \%LS)$$

$$\%data = \text{Percent of Data Accesses} = \%LS / (1 + \%LS)$$

$$\%LS = \text{Frequency of Load and Store instructions}$$

- **每条指令的L1 失效次数:**

$$\text{Misses per Instruction}_{L1} = \text{Miss Rate}_{L1} \times (1 + \%LS)$$

$$\text{Misses per Instruction}_{L1} = \text{Miss Rate}_{I\text{-Cache}} + \%LS \times \text{Miss Rate}_{D\text{-Cache}}$$



具有二级Cache的AMAT举例

- **Problem: 计算AMAT**

- I-Cache 失效率 = 1%, D-Cache失效率 = 10%
- L2 Cache失效率 = 40%
- L1 命中时间 = 1 cycle (I-Cache 和D-Cache相同)
- L2 命中时间 = 8 cycles, L2 失效开销 = 100 cycles
- Load + Store 指令频度 = 25%

- **Solution:**

- 平均每条指令访存次数 = $1 + 25\% = 1.25$
- 平均每条指令的失效次数 = $1\% + 25\% \times 10\% = 0.035$
- L1的失效率 = $0.035 / 1.25 = 0.028$
- L1的失效开销 = $8 + 0.4 \times 100 = 48$ cycles
- AMAT = $1 + 0.028 \times 48 = 2.344$



Memory Stall Cycles Per Instruction

- **Memory Stall Cycles per Instruction**

= Memory Access per Instruction \times Miss RateL1 \times Miss PenaltyL1

= $(1 + \%LS) \times$ Miss RateL1 \times Miss PenaltyL1

= $(1 + \%LS) \times$ Miss RateL1 \times (Hit TimeL2 + Miss RateL2 \times Miss PenaltyL2)

- **Memory Stall Cycles per Instruction**

= Misses per InstructionL1 \times Hit TimeL2 +

Misses per InstructionL2 \times Miss PenaltyL2

Misses per InstructionL1 = $(1 + \%LS) \times$ Miss RateL1

Misses per InstructionL2 = $(1 + \%LS) \times$ Miss RateL1 \times Miss RateL2

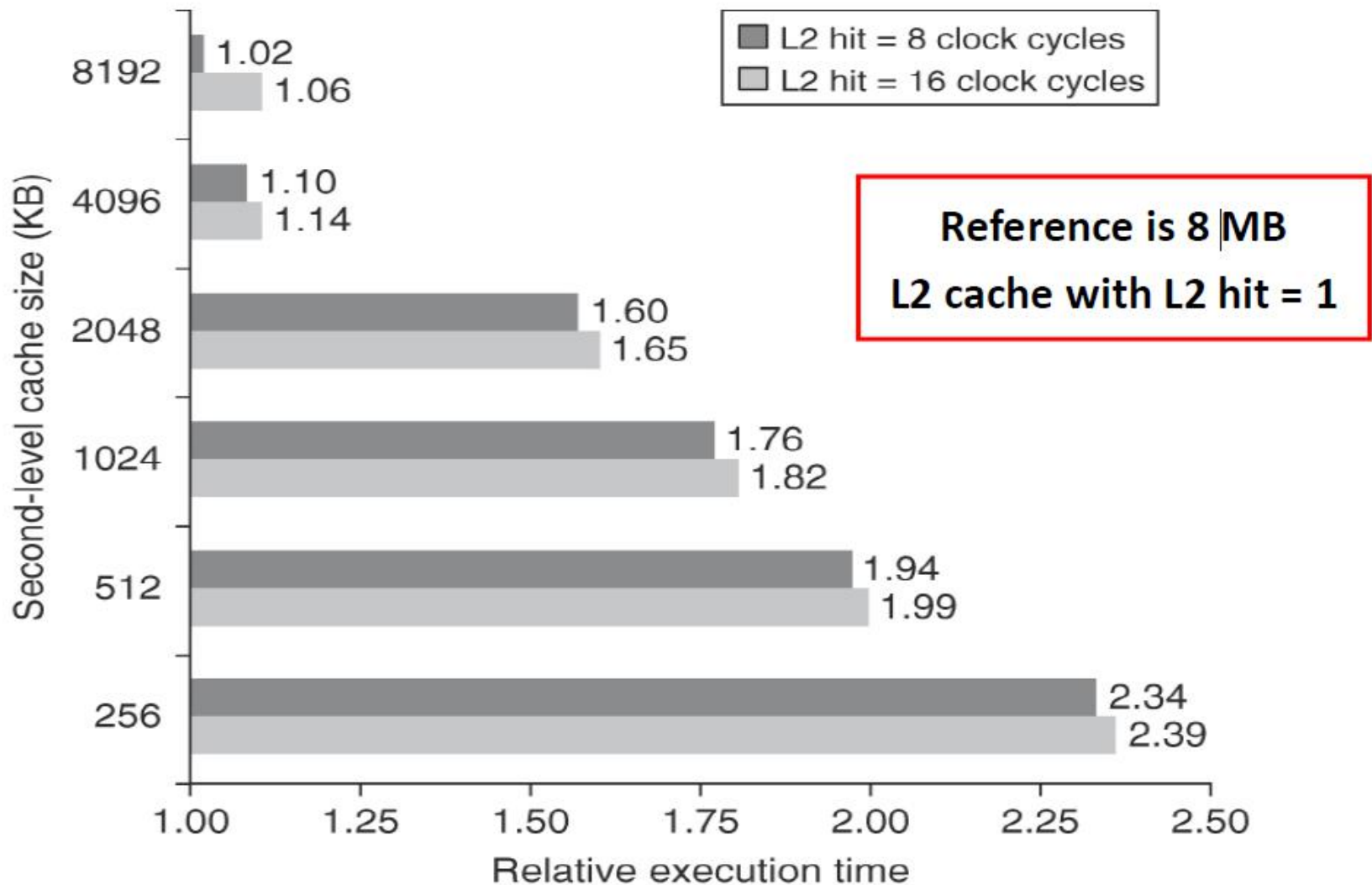


两级Cache的性能

- **Problem: 程序运行产生1000个存储器访问**
 - I-Cache misses = 5, D-Cache misses = 35, L2 Cache misses = 8
 - L1 Hit = 1 cycle, L2 Hit = 8 cycles, L2 Miss penalty = 80 cycles
 - Load + Store frequency = 25%, CPI_{execution} = 1.1 (perfect cache)
 - 计算memory stall cycles per instruction 和有效的CPI
 - 如果没有L2 cache, 有效的CPI是多少?
- **Solution:**
 - L1 Miss Rate = $(5 + 35) / 1000 = 0.04$ (or 4% per access)
 - L1 misses per Instruction = $0.04 \times (1 + 0.25) = 0.05$
 - L2 misses per Instruction = $(8 / 1000) \times 1.25 = 0.01$
 - Memory stall cycles per Instruction = $0.05 \times 8 + 0.01 \times 80 = 1.2$
 - CPI_{L1+L2} = $1.1 + 1.2 = 2.3$, CPI/CPI_{execution} = $2.3/1.1 = 2.1x$ slower
 - CPI_{L1only} = $1.1 + 0.05 \times 80 = 5.1$ (worse)



不同大小的L2cache对应的执行时间

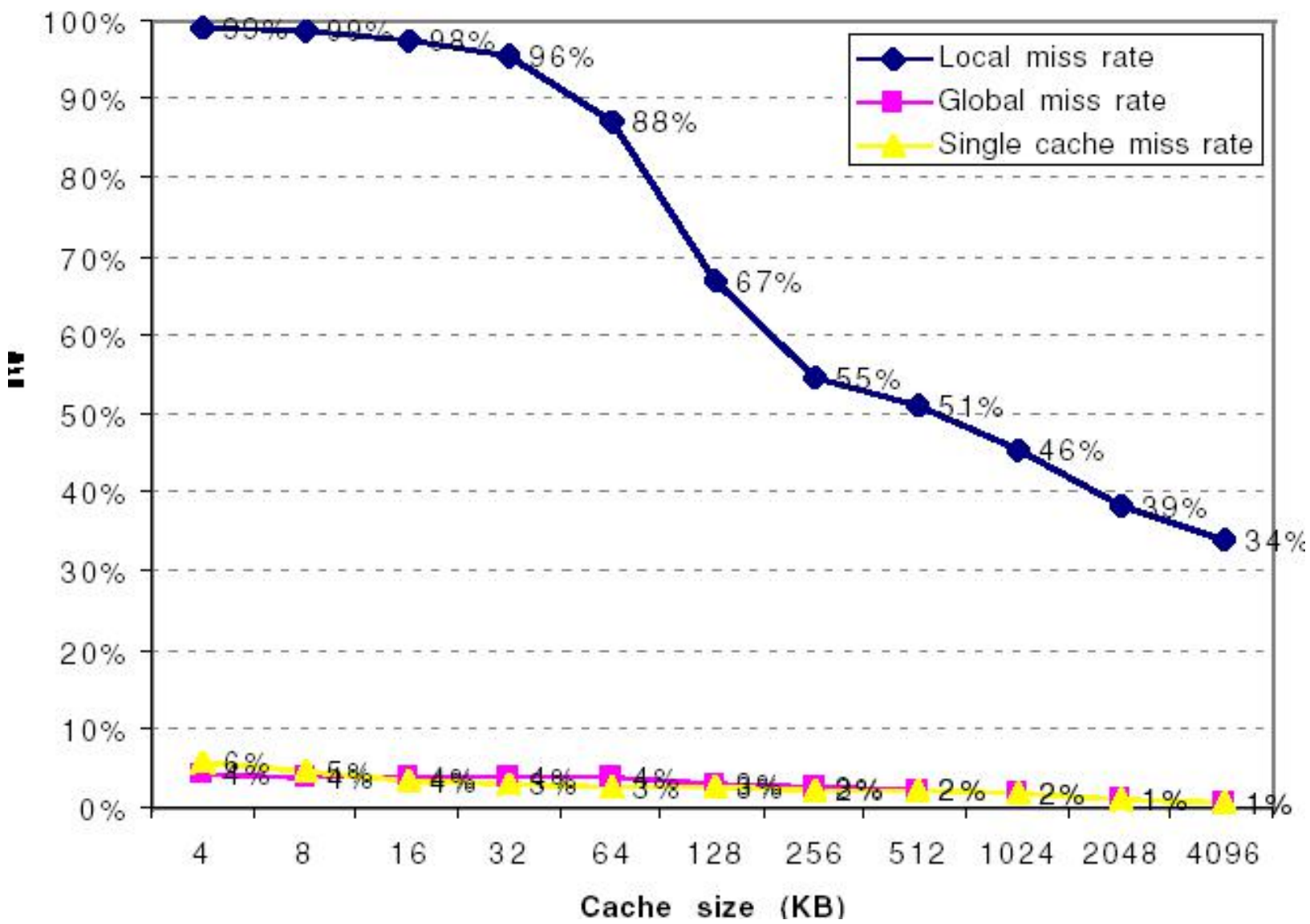


Copyright © 2012, Elsevier Inc. All rights reserved.

计算机体系结构



Miss rates versus cache size for multilevel caches





两级Cache的一些研究结论

- **在L2比L1大得多得情况下，两级Cache全局失效率 和容量与第二级Cache相同的单级Cache的失效率接近**
- **局部失效率不是衡量第二级Cache的好指标**
 - 它是第一级Cache失效率的函数
 - 不能全面反映两级Cache体系的性能
- **第二级Cache设计需考虑的问题**
 - **容量**：一般很大，可能没有容量失效，只有强制性失效和冲突失效
 - 相联度对第二级Cache的作用
 - Cache可以较大，以减少失效次数
 - **多级包容性问题**：第一级Cache中的数据是否总是同时存在于第二级Cache中。
 - 如果L1和L2的块大小不同，增加了多级包容性实现的复杂性



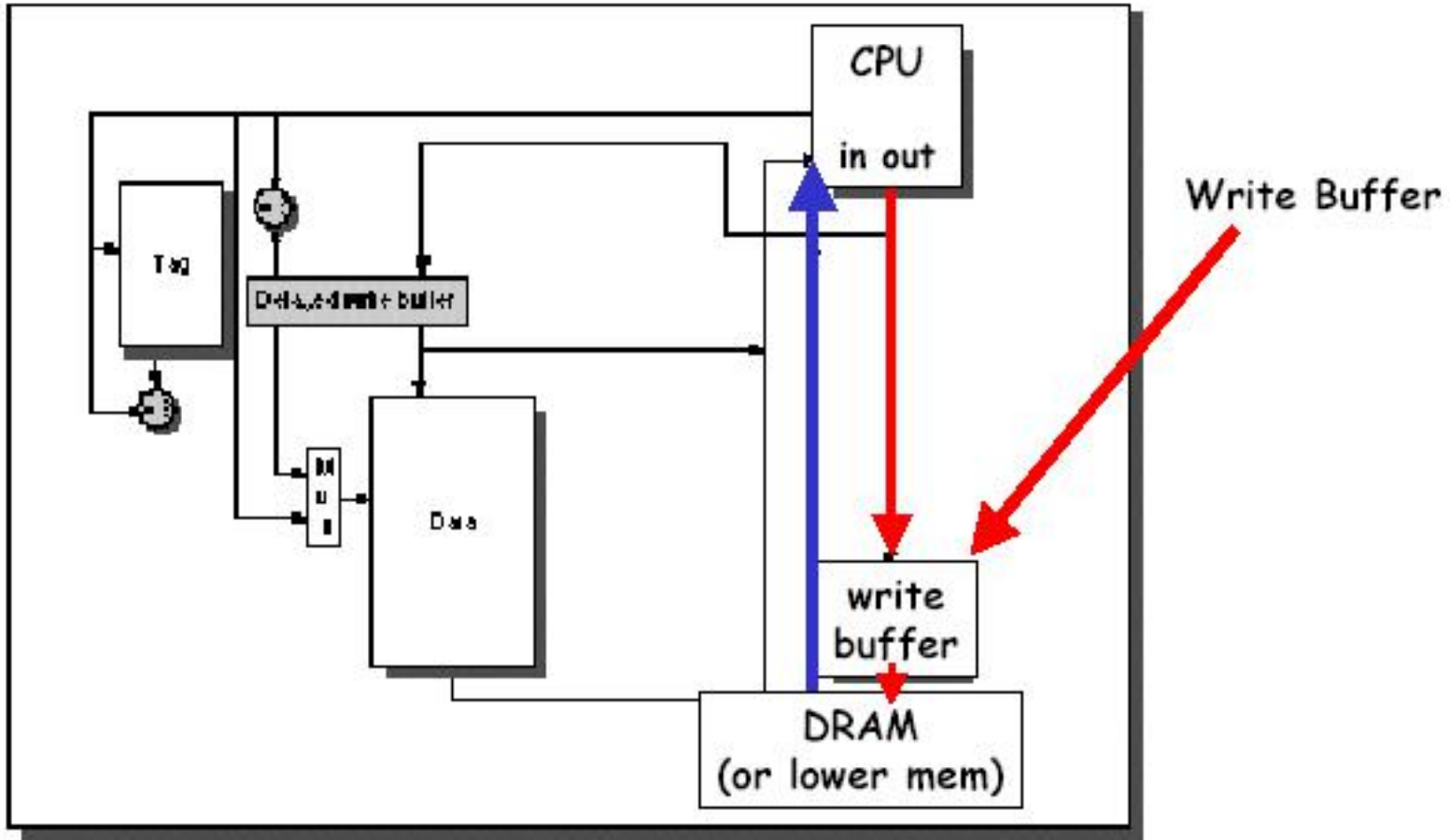
多级Cache举例

Given the data below, what is the impact of second-level cache associativity on its miss penalty?

- Hit timeL2 for direct mapped = 10 clock cycles
 - Two-way set associativity increases hit time by 0.1 clock cycles to 10.1clock cycles
 - Local miss rateL2 for direct mapped = 25%
 - Local miss rateL2 for two-way set associative = 20%
 - Miss penaltyL2 = 100 clock cycles
-
- **结论：提高相联度，可减少第一级Cache的失效开销**
 - **第二级Cache特点：容量大，高相联度，块较大，重点减少失效次数。**



让读优先于写图示





让读失效优先于写

- **由于读操作为大概率事件，需要读失效优先，以提高性能**
- **Write-Through Cache -> Write Buffer (写缓冲)，特别对写直达法更有效**
 - Write Buffer: CPU不必等待写操作完成，即将要写的数据和地址送到Write Buffer后，CPU继续作其他操作。
 - 写缓冲导致对存储器访问的复杂化
 - 在读失效时写缓冲中可能保存有所读单元的最新值，还没有写回
 - 例如，直接映射、写直达、512和1024映射到同一块。则
 - SW R3, 512(R0)
 - LW R1, 1024(R0) 失效
 - LW R2, 512(R0) 失效
 - 解决问题的方法
 - 推迟对读失效的处理，直到写缓冲器清空，导致新的问题——读失效开销增大。
 - 在读失效时，检查写缓冲的内容，如果没有冲突，而且存储器可访问，就可以继续处理读失效
 - 写回法时，也可以利用写缓冲器来提高性能
 - 把脏块放入缓冲区，然后读存储器，最后写存储器
- **Write-Back Cache -> Victim Buffer**
 - 被替换的脏块放到了victim buffer
 - 在脏块被写回前，需要处理读失效
 - 问题: victim buffer 可能含有该读失效要读取的块
 - Solution: 查找victim buffer，如果命中直接将该块调入Cache



4.2 Cache 基本优化方法

Cache
性能分析

降低
失效率

减少
失效开销

缩短
命中时间

平均访存时间 = **命中时间** + 失效率 × 失效开销



缩短命中时间

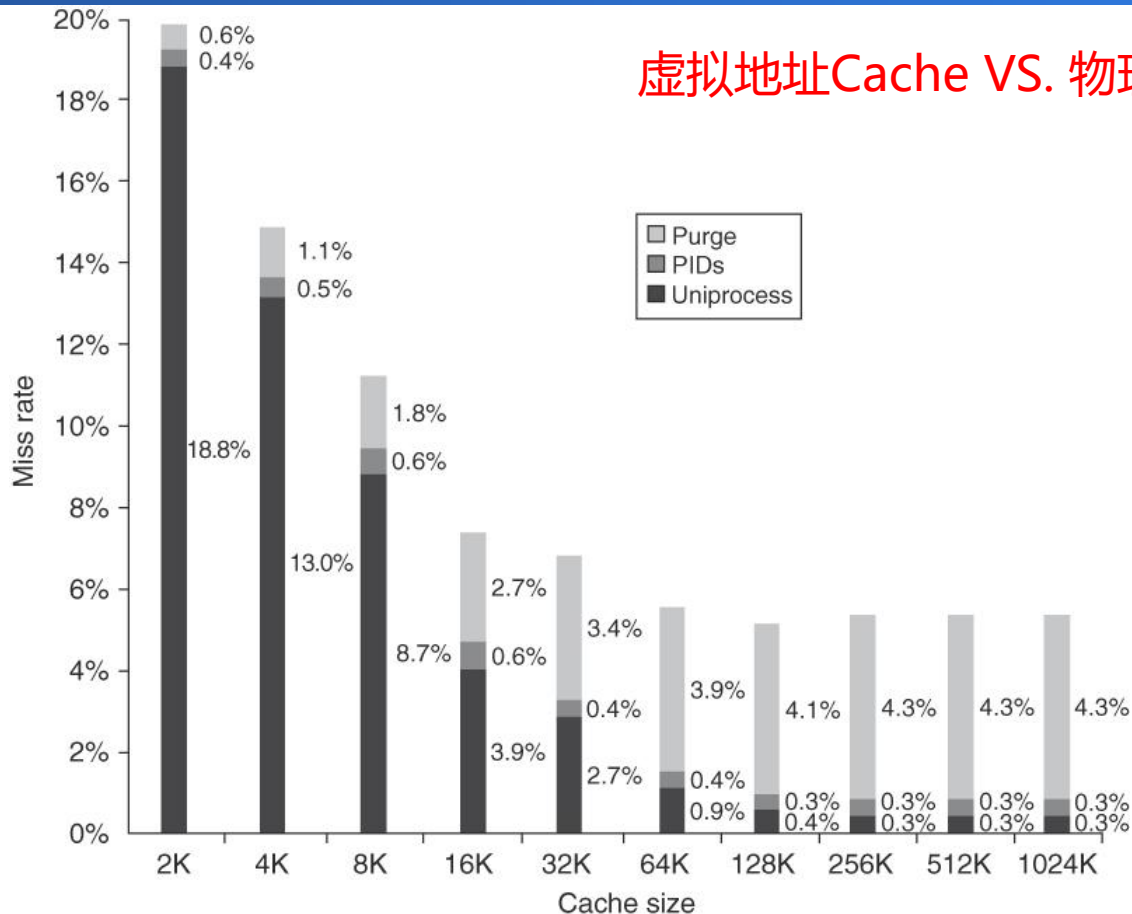
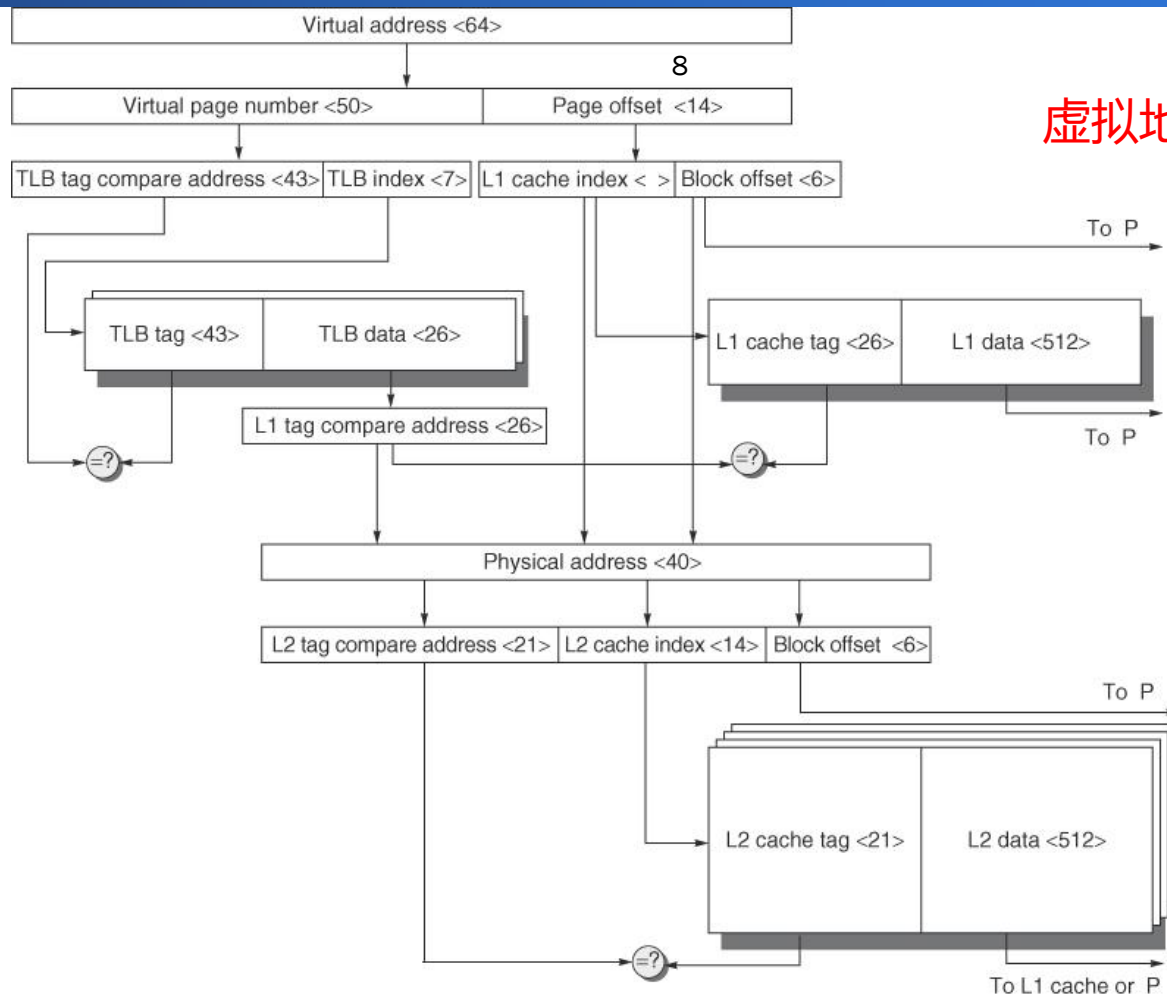


Figure B.16 Miss rate versus virtually addressed cache size of a program measured three ways: without process switches (uniprocess), with process switches using a process-identifier tag (PID), and with process switches but without PIDs (purge). PIDs increase the uniprocess absolute miss rate by 0.3% to 0.6% and save 0.6% to 4.3% over purging. Agarwal [1987] collected these statistics for the Ultrix operating system running on a VAX, assuming direct-mapped caches with a block size of 16 bytes. Note that the miss rate goes up from 128K to 256K. Such nonintuitive behavior can occur in caches because changing size changes the mapping of memory blocks onto cache blocks, which can change the conflict miss rate.



虚拟地址转换与Cache定位 并行

Figure B.17 The overall picture of a hypothetical memory hierarchy going from virtual address to L2 cache access. The page size is 16 KB. The TLB is two-way set associative with 256 entries. The L1 cache is a direct-mapped 16 KB, and the L2 cache is a four-way set associative with a total of 4 MB. Both use 64-byte blocks. The virtual address is 64 bits and the physical address is 40 bits.



小结：Cache 性能分析

- **CPU time = (CPU execution clock cycles + Memory stall clock cycles) x clock cycle time**
- **Memory stall clock cycles = (Reads x Read miss rate x Read miss penalty + Writes x Write miss rate x Write miss penalty)**
- **Memory stall clock cycles = Memory accesses x Miss rate x Miss penalty**
- **Different measure: AMAT**

Average Memory Access time (AMAT) = Hit Time + (Miss Rate x Miss Penalty)

- **Note: *memory hit time is included in execution cycles.***



小结：基本Cache优化方法

• 基本Cache优化方法

降低失效率：引起失效的3C

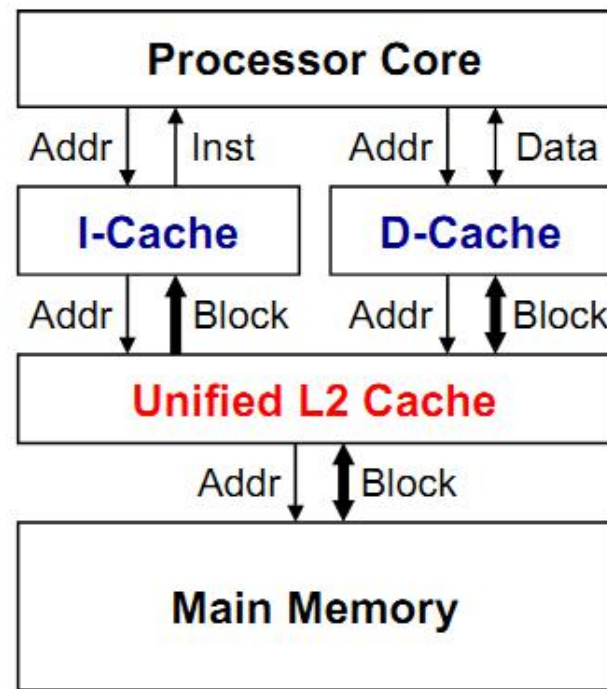
- 1、增加Cache块的大小
- 2、增大Cache容量
- 3、提高相联度

减少失效开销

- 4、多级Cache
- 5、使读失效优先于写失效

缩短命中时间

- 6、避免在索引缓存期间进行地址转换





Acknowledgements

- **These slides contain material developed and copyright by:**
 - John Kubiatowicz (UCB)
 - Krste Asanovic (UCB)
 - David Patterson (UCB)
 - Chenxi Zhang (Tongji)
- **UCB material derived from course CS152、CS252、CS61C**
- **KFUPM material derived from course COE501、COE502**