

# Impact on Performance

- 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
- Suppose that 10% of memory operations get 50 cycle miss penalty
- Suppose that 1% of instructions get same miss penalty
- $CPI = \text{ideal CPI} + \text{average stalls per instruction}$ 

$$1.1(\text{cycles/ins}) + [0.30(\text{DataMops/ins}) \times 0.10(\text{miss/DataMop}) \times 50(\text{cycle/miss}) + 1(\text{InstMop/ins}) \times 0.01(\text{miss/InstMop}) \times 50(\text{cycle/miss})]$$

$$= (1.1 + 1.5 + .5) \text{ cycle/ins} = 3.1$$

每条指令都要取指  
Memory 操作  
Cache = Miss

取instr load/store 共1.3  
有1.3个的指令

会用到取instr  
0.3的指令load/store  
共1.3的指令

Van Neumann

① 层次内存  
会按需求

② 任何都是data

58% of the time the proc is stalled waiting for memory!

$$AMAT = (1/1.3) \times [1 + 0.01 \times 50] + (0.3/1.3) \times [1 + 0.1 \times 50] = 2.54$$

1/24/01

U.S.

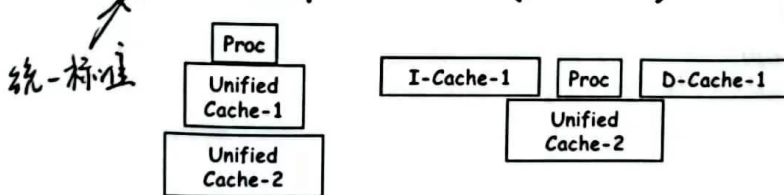
Arken

所有内存不超过1KB

## Example: Harvard Architecture

instr. ≠ data.

- Unified vs Separate I&D (Harvard)



- Table on page 384:
  - 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)

$$AMAT_{\text{Harvard}} = 75\% \times (1 + 0.64\% \times 50) + 25\% \times (1 + 6.47\% \times 50) = 2.05$$

$$AMAT_{\text{Unified}} = 75\% \times (1 + 1.99\% \times 50) + 25\% \times (1 + 1.99\% \times 50) = 2.24$$

1/24/01

CS252/Kubiatowicz Lec 3.10

单端口：同时读指令和数据会出问题

SMC: self Modify co  
超级扇区防篡改

$$A.M.A.T = \underbrace{T_{hit}}_{\text{Direct Mapping}} + \underbrace{1}_{\text{Associative}} \cdot \underbrace{T_{penalty}}_{\text{Memory}}$$

## Reducing Misses

### • Classifying Misses: 3 Cs

- **Compulsory**—The first access to a block is not in the cache, so the block must be brought into the cache. Also called cold start misses or first reference misses.

(Misses in even an Infinite Cache)

- **Capacity**—If the cache cannot contain all the blocks needed during execution of a program, capacity misses will occur due to blocks being discarded and later retrieved.

(Misses in Fully Associative Size X Cache)

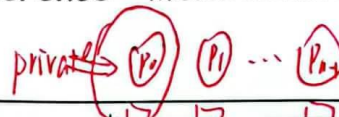
- **Conflict**—If block-placement strategy is set associative or direct mapped, conflict misses (in addition to compulsory & capacity misses) will occur because a block can be discarded and later retrieved if too many blocks map to its set. Also called collision misses or interference misses.

(Misses in N-way Associative, Size X Cache)

### • More recent, 4th "C":

- **Coherence** - Misses caused by cache coherence.

1/24/01



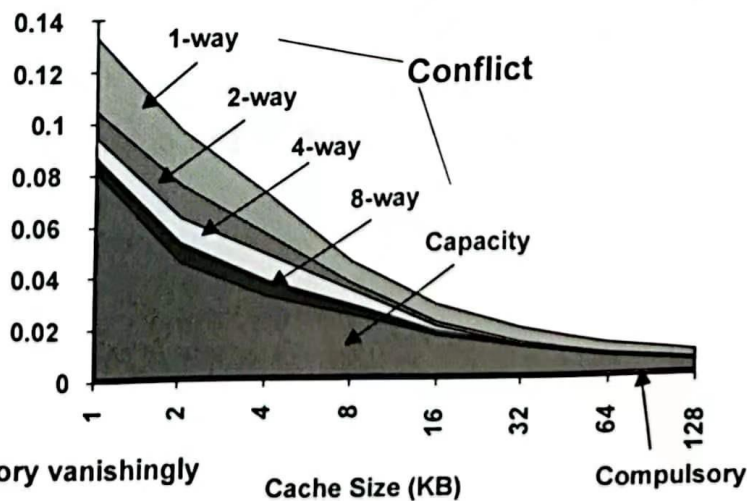
连接不同程序的  
cache之间有不同

CS252/Kubiatowicz  
Lec 3.14

not coherent

Memory public & shared

## 3Cs Absolute Miss Rate (SPEC92)

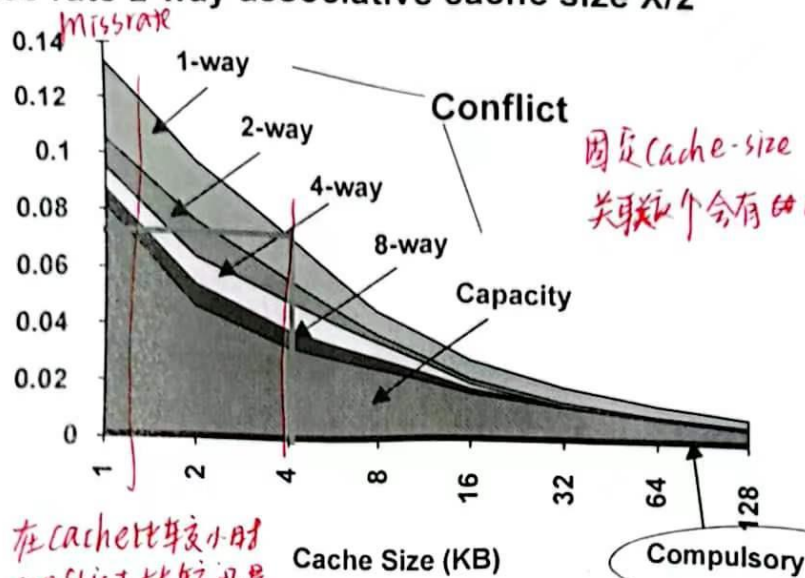


1/24/01

CS252/Kubiatowicz  
Lec 3.15

## 2:1 Cache Rule

miss rate 1-way associative cache size X  
= miss rate 2-way associative cache size X/2



固定 cache-size  
关联数 ↑ 会有 miss-rate ↓

在 cache 比较小时  
conflict 比较明显  
因此提高关联数 有效  
解决 conflict 的问题

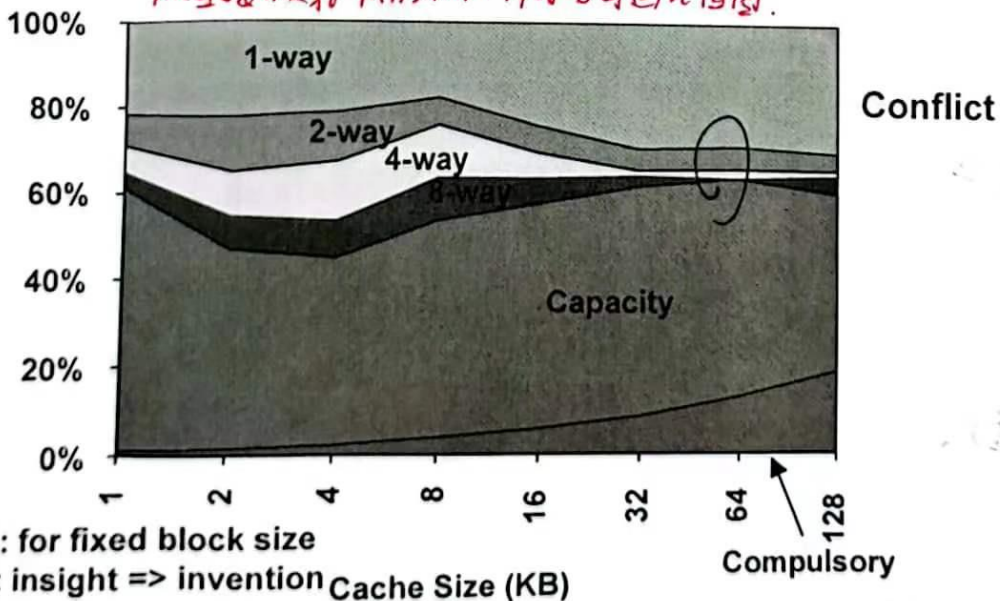
block 块大小不变, 第一次  
miss 率不变

1/24/01

CS252/Kubiatowicz  
Lec 3.16

## 3Cs Relative Miss Rate

以直接映射 miss rate 作为分母画比值图



Flaws: for fixed block size

Good: insight => invention

1/24/01

CS252/Kubiatowicz  
Lec 3.17



## How Can Reduce Misses?

- 3 Cs: Compulsory, Capacity, Conflict
- In all cases, assume total cache size not changed:
- What happens if:
  - 1) Change Block Size:  
Which of 3Cs is obviously affected?
  - 2) Change Associativity:  
Which of 3Cs is obviously affected?
  - 3) Change Compiler:  
Which of 3Cs is obviously affected?

1/24/01

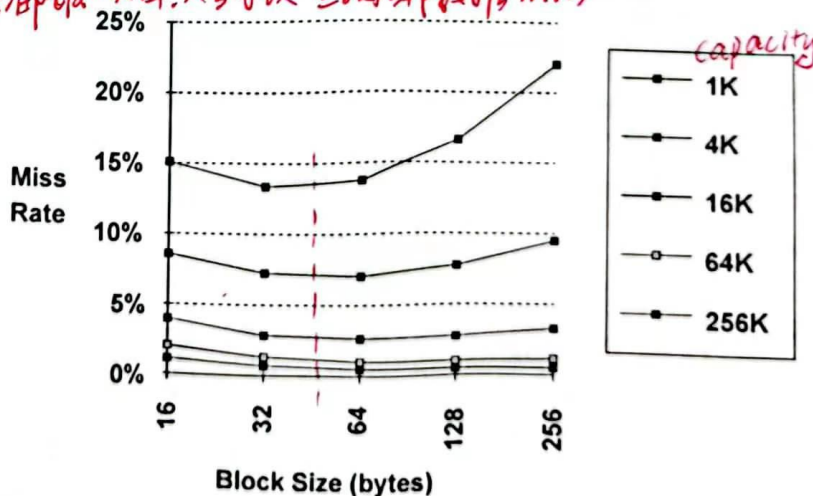
CS252/Kubiatowicz  
Lec 3.18

block 有一个 tag, 标志着整块的信息. 一个 block 中的数据在 memory 中始终是连续的

block 小的时候, 受 capacity 限制  
block 略大, 可以发挥 locality, 拿进来更多附近数据

### 1. Reduce Misses via Larger Block Size

但 block 过大, 每 miss 一个数据  
整块 block 都被换掉, 又会导致一些捆绑数据 miss



← 破坏空间本地性

→ :  
捆绑一起  
一个 miss 换掉整块

1/24/01

CS252/Kubiatowicz  
Lec 3.19

$$\frac{1}{2^{\log}} \rightarrow \frac{2}{2^{\log+1}}$$

由本地性等数据特点

## 2. Reduce Misses via Higher Associativity

### • 2:1 Cache Rule:

- Miss Rate DM cache size N - Miss Rate 2-way cache size N/2

### • Beware: Execution time is only final measure!

- Will Clock Cycle time increase?
- Hill [1988] suggested hit time for 2-way vs. 1-way external cache +10%, internal + 2%

Miss Rate: 并不是冷冰冰的概率比值  
(由于数据使用的特点和分布)  $\Leftrightarrow$  本地性

1/24/01

CS252/Kubiatowicz  
Lec 3.20

随着路数增加, missrate 会 ↓

$$A.M.A.T. = \bar{T}_{hit} + \underbrace{\eta_{missrate}}_{\text{随着路数增加, 找到 hit 的时间}} \cdot T_{penalty}$$

相互作用  
先 ↑ 后 ↓

## Example: Avg. Memory Access Time vs. Miss Rate

- Example: assume CCT = 1.10 for 2-way, 1.12 for 4-way, 1.14 for 8-way vs. CCT direct mapped

Cache Size (KB)	Associativity			
	1-way	2-way	4-way	8-way
1	2.33	2.15	2.07	2.01
2	1.98	1.86	1.76	1.68
4	1.72	1.67	1.61	1.53
8	1.46	1.48	1.47	1.43
16	1.29	1.32	1.32	1.32
32	1.20	1.24	1.25	1.27
64	1.14	1.20	1.21	1.23
128	1.10	1.17	1.18	1.20

单位 cycle

(Red means A.M.A.T. not improved by more associativity)

1/24/01

CS252/Kubiatowicz  
Lec 3.21

compulsory 跟 block 大小有关  
conflict 跟 关联度 有关  
capacity 跟 总容量 有关



Write back & Write through.

如果某 cache 的数据  $x \rightarrow x'$   
而 memory 还未被修改, 这比 cache  
称为 dirty cache.

dirty block (write through).

### 3. Reducing Misses via a "Victim Cache"

- How to combine fast hit time of direct mapped yet still avoid conflict misses?
- Add buffer to place data discarded from cache
- Jouppi [1990]: 4-entry victim cache removed 20% to 95% of conflicts for a 4 KB direct mapped data cache
- Used in Alpha, HP machines

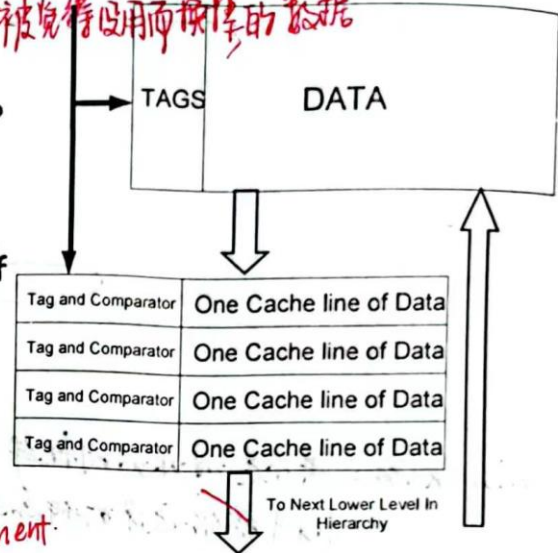
在工程实现中实现比较相似

Victim Cache vs. Write Buffer

victim cache 是 cache 满之后  
被替换掉的数据暂存到 victim  
cache 中, 如果下次再用可以拿回

write buffer 用于 write through.  
在 cache 中修改, 肯定得同步到  
memory 中. 先放到 buffer 中  
找合适的位置放回.

被替换使用而替换的数据

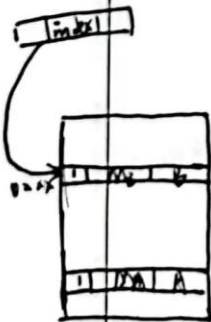


victim cache & memory  
脏数据不是必须的, victim cache

相当于被暂时存放 cache 废弃数据  
的垃圾桶, 而 write buffer 要 & me  
相连.

### 4. Reducing Misses via "Pseudo-Associativity"

- How to combine fast hit time of Direct Mapped and have the lower conflict misses of 2-way SA cache?
- Divide cache: on a miss, check other half of cache to see if there, if so have a pseudo-hit (slow hit)



对全关联的一个改进, 利用空位置  
Hit Time 如果出现 conflict 不先 replace.  
Pseudo Hit Time 在这之前翻转 tag 的首位用另一半空间  
Miss Penalty

Drawback: CPU pipeline is hard if hit takes 1 or 2 cycles  
- Better for caches not tied directly to processor (L2)  
- Used in MIPS R1000 L2 cache, similar in UltraSPARC

保持直接映射高速的同时  
也利用另一半减少 Conflicts Miss

增加一个 index MSB  
防止出现问题.  
防止被占用的 (转就以 0/1 开头  
为要用 cache)

Locality

提高了空间本地性 提高了 compulsory

把相近的数据拿进来

stream buffer

CPU

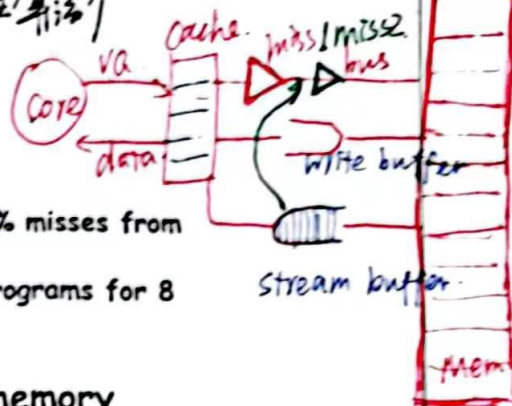
compulsory is miss rate ↓.

## 5. Reducing Misses by Hardware Prefetching of Instructions & Datals

### • E.g., Instruction Prefetching

线性算法 ↑

- Alpha 21064 fetches 2 blocks on a miss
- Extra block placed in "stream buffer"
- On miss check stream buffer



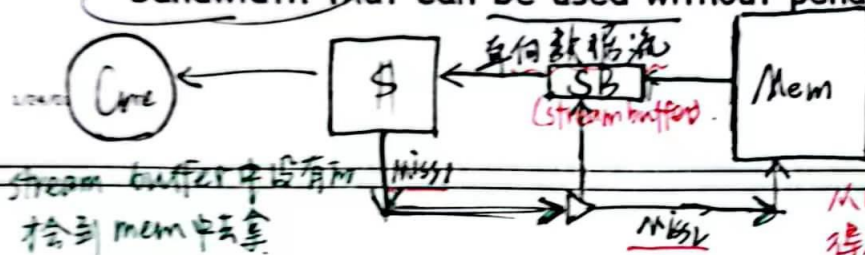
### • Works with data blocks too:

cache block 有 4K  
block 的 stream buffer  
也要 4 个 block

Jouppi [1990] 1 data stream buffer got 25% misses from 4KB cache; 4 streams got 43%

Palacharla & Kessler [1994] for scientific programs for 8 streams got 50% to 70% of misses from 2 64KB, 4-way set associative caches

### • Prefetching relies on having extra memory bandwidth that can be used without penalty



CS252/Kubiatowicz  
Lec 3.24

stream buffer 中只有刚刚

从 mem 拿过来的数据，还没来得及放入 cache，在 cache 中 miss 1.

Compiler

就寄去 stream buffer

在 stream buffer 中 miss 2 才会去 mem

## 6. Reducing Misses by Software Prefetching Data

### • Data Prefetch

- Load data into register (HP PA-RISC loads)
- Cache Prefetch: load into cache (MIPS IV, PowerPC, SPARC v. 9)
- Special prefetching instructions cannot cause faults; a form of speculative execution

### • Prefetching comes in two flavors:

Nothing!

- Binding prefetch: Requests load directly into register
- » Must be correct address and register!
- Non-Binding prefetch: Load into cache.
- » Can be incorrect. Frees HW/SW to guess!

LD R<sub>i</sub>, R<sub>j</sub>(#n)

LD R<sub>j</sub>(#n)

### • Issuing Prefetch Instructions takes time

- Is cost of prefetch issues < savings in reduced misses?
- Higher superscalar reduces difficulty of issue bandwidth