

5.16

6 因为使用中间位作为组索引可以使相邻的地址映射到同一组中，利用局部性原理提高缓存效率  
高位作为标签可以在多级缓存系统中，将低层缓存中相同地址的块合并，提高缓存灵活性

7 这样做的好处在于：1°可以简化虚拟内存到物理内存的地址转换 2°可以用虚拟地址的  
页偏移位数来访问缓存，提高缓存的访问速度 3°使每个缓存块可以在储一个完整的  
虚拟页面，提高缓存利用率。

8.(1) 平均访问延时为  $T = 1 \times 97\% + 1/10 \times 3\% = 4.27$  周期

(2)  $1GB = 2^{30} KB$   $\therefore$  有  $\frac{2^30}{2^{10}} = 0.00006$  的几率命中

∴ 平均访问延时为  $T = 1 \times 0.00006 + 1/10 \times 0.99994 = 109.9934$  周期

(3) 局部性原理是实际处理器能通过缓存极大减少存储系统平均访问延时的根本原因。因为  
时间空间的局部性原理，从 cache 的命中率才可能达到 97%，否则只有  $\frac{1}{2}$

(4) 设命中率为  $\alpha$   $1 \cdot \alpha + (1/10 \cdot (1-\alpha)) \leq 1/5 \quad \alpha \geq \frac{5}{109} \approx 0.0459$

9 地址位数 (bit)	缓存大小 (KB)	块大小 (Byte)	相联度	组数量	组索引位数	标签位数	偏移位数
32	4	64	2	32	5	21	6
32	4	64	8	8	3	23	6
32	4	64	全相联	1	10	26	6
32	16	64	1	256	8	18	6
32	16	128	2	64	6	19	7
32	64	64	4	16	8	18	6
32	64	64	16	4	6	20	6
32	64	128	16	2	5	20	7

$$(1) 0.22(1-p_1) + \frac{(0.22+0.22)p_1}{1-p_1} < 0.52(1-p_2) + \frac{(0.52+0.52)p_2}{1-p_2}$$

$$0.22 - 0.22p_1 + 0.22p_1 < 0.52 - 0.52p_2 + 0.52p_2$$

$$0.22 < 0.52, p_1 - p_2 < 0.003$$

$$(2) 0.22(1-p_1) + 0.22(k+1)p_1 < 0.52(1-p_2) + 0.52(k+1)p_2$$

$$0.22kp_1 - 0.52kp_2 < 0.3$$

11  $0x$  为 16 位制 共 16 个块

直接映射时，有 16 个组  $0x100$  放入组 1  $0x1005$  放入组 5  $0x104$  替换放入组 1， $0x1045$ ，  
 $0x1005$   $0x2005$   $0x1005$  依次替换组 5  $\therefore$  共 5 次替换

2 路组相联 8 个组  $0x100$  放入组 1  $0x1005$  放入组 5  $0x104$  放入组 1  $0x1045$  放入组 5  
 后面 3 个请求都会发生替换  $\therefore$  共 3 次替换

4 路组相联 4 个组 mod 4  $0x100$ ,  $0x1005$ ,  $0x104$ ,  $0x1045$  都放入组 1, 后面 3 个请求也要放入组 1  $\therefore$  3 次替换

8 路组相联 8 个组 mod 2 除，所有数据都要放到组 1，组 1 有 8 个位置，可以放下  
 $\therefore$  0 次替换

12 块 16 Byte, 总容量 256 Byte  $\Rightarrow$  16 个块 int 32-bit 54 Byte/e 每块可放 4 个 int 32-bit

A: 8 个组 每路可放  $8 \times 4 = 32$  个 int 共 96 个 int，将其三等份，称为 I, II, III

A 中 第一路存 I, 第二路存 II, 存 III 时根据 LRU, 放在第 I 路。之后第 II 路存 II, 第 III 路存 III

$\begin{matrix} 1 & 2 \\ \text{I} & \text{II} \end{matrix} \rightarrow \begin{matrix} 1 & 2 \\ \text{II} & \text{III} \end{matrix} \rightarrow \begin{matrix} 1 & 2 \\ \text{III} & \text{I} \end{matrix} \rightarrow \begin{matrix} 1 & 2 \\ \text{II} & \text{I} \end{matrix}$   $N=100$  时，缺失率  $\approx \frac{1}{4}$

B: 16 个组  $\begin{matrix} 1 \\ \text{I} \\ \text{II} \end{matrix} \rightarrow \begin{matrix} 1 \\ \text{II} \\ \text{III} \end{matrix} \rightarrow \begin{matrix} 1 \\ \text{III} \\ \text{I} \end{matrix} \rightarrow \begin{matrix} 1 \\ \text{I} \\ \text{II} \end{matrix}$  缺失率  $\approx \frac{1}{4}$

13 改为  $\text{for}(\text{int } j=0; j < 128; j ++)$

$\text{for}(\text{int } i=0; i < 64; i ++)$

$$A[j][i] = A[j][i] + 1;$$

14 块32字节  $\Rightarrow$  每块8个int 4KB中共128个块

(1) 优化前 直接映射 内存中  $A[0][0] \sim A[0][B]$  占  $64 \times 4 = 2^8$  Byte

$\therefore 2^{12} \div 2^8 = 16$ , 每有块地址  $j \pmod{16}$  时会存放在相同块

$\because j=0 \sim 127$ ,  $i=0 \sim 64$   $\therefore$  每次都缺失, 次数为  $64 \times 128 = 8192$  次

优化后 每块8个int 1. 缺失率  $\frac{1}{8}$  2. 缺失次数  $64 \times 128 \div 8 = 1024$  次

(2) 全相联 优化前 缺失次数为  $64 \times 128 \div 8 = 1024$  次

优化后 同上, 为  $64 \times 128 \div 8 = 1024$  次

(3) 优化前: 要扩容到  $4KB \times \frac{128}{16} = 32KB$  优化后仍为4 KB

input                          output

15 列0 列1 列2 列3      列0 列1 列2 列3

行0 miss hit miss hit      miss miss miss miss

行1 miss hit miss hit      miss miss miss miss

行2 miss hit miss hit      miss miss miss miss

行3 miss hit miss hit      miss miss miss miss

16(a) 4个hit/块 共32个块 3路组相联, 共16组

缺失率  $\frac{1}{4}$ , 命中率  $\frac{3}{4}$

(2) 不能, 因为增加缓存总大小无法减少强制缺失

(3) 可以, 块变大可以减少强制缺失