

6. 因为组索引越靠前的位数变化的范围越小，可以更好地利用局部性原理，提高缓存的命中率。而标签则需要存储的信息较多，需要更多位数来表示。

7. 使地址的组索引和块内偏移的总位数与虚拟内存系统的页偏移位数相同，可以使缓存系统的映射方式与虚拟内存系统的映射方式相同，方便管理与协调。

8) 缓存平均访问延时 = 缓存命中时间 + 缓存缺失时间 × 缓存缺失率 = 1 + 10 × 3% + 105 × (1-3%) = 1.05 个周期

2) 每个块 64 字节，64KB 的 L1 缓存可以缓存 1024 块，可以覆盖整个数组。由于完全随机访问，每个块的访问概率相等，因此平均访问延时 = 1 + 105 × 3% = 4.15 个周期

3) 局部性原理好坏直接影响了缓存的命中率和缺失率，进而影响处理器访存的性能。当访问随机性很高时，缓存的效果就不如预期，当程序的访问具有局部性时，缓存会取得更好的效果。

4) 程序的平均缓存命中率需要高于 3%，才能让存储系统在使用 L1 时获得性能收益。

|   | 组数量 | 组索引位数 Bit | 行地址位数 Bit | 偏移位数 Bit |
|---|-----|-----------|-----------|----------|
| 1 | 32  | 5         | 21        | 6        |
| 2 | 8   | 3         | 23        | 6        |
| 3 | 1   | 0         | 26        | 6        |
| 4 | 256 | 8         | 18        | 6        |
| 5 | 64  | 6         | 19        | 7        |
| 6 | 256 | 8         | 18        | 6        |
| 7 | 64  | 6         | 20        | 6        |
| 8 | 32  | 5         | 20        | 7        |

10) (1) 系统 A 的平均内存访问时间 = 缓存缺失时间 × 缓存缺失率 + 缓存缺失代价 × 缓存缺失率 + 缓存命中时间 = 0.22 + 100P1 + 110P2 → 0.52 + 100P2, 即  $P1 > 5P2 - 1.1$   
(2) 系统 A 的平均内存访问时间 = 缓存命中时间 + 缓存缺失时间 × 缓存缺失率 + 缓存缺失代价 × 缓存缺失率 × k = 0.22 + 0.52k + 100P1k + 110P2k, 而系统 B 的平均内存访问时间为  $0.52 + 0.52k + 100P2k$ , 若两者相等，则  $P1 > 1.75P2 - 0.3$

11. 直接映射，缓存容量为  $16 \times 64 \div 8 = 128$  字节。  
2 路组相联，缓存容量为  $2 \times 16 \times 64 \div 8 = 256$  字节，块替换次数为 1。  
4 路组相联，缓存容量为  $4 \times 16 \times 64 \div 8 = 512$  字节，块替换次数为 0。  
同理 8 路组相联，缓存容量为 1024 字节，块替换次数为 0。

12. 缓存 A 为 2 路组相联，总共可存放 8 块，缓存 B 为直接映射，总共可存放 16 块。  
优化前，数组 A 在缓存 A 中的缺失次数为 640，缓存 B 中的缺失次数为 1600。  
优化后，数组 A 在缓存 A 中的缺失次数为 128，缓存 B 中的缺失次数为 2510。

13. 将内层循环中访问 A 数组的方式修改为按照行优先，即将  $A[i][j]$  修改为  $A[i][j]$ ，可以更好地利用行主元局部性，减少缓存缺失次数。

14. (1) 优化前，总共发生 256 次缓存缺失；  
优化后，缓存缺失次数减少到 128 次，因为每次循环只需要缓存一行或一列 (16 字节)，因为块大小为 32 字节，所以每个块可以缓存 2 个元素，故需 128 个块。

(2) 优化前，总共发生 256 次缓存缺失；  
优化后，缓存缺失次数减少到 0 次；因为全相联缓存可以缓存所有数据，不会发生替换。

(3) 对于块大小 32 字节的直接映射缓存，每个块可以包含 2 个元素；缓存所需块数为  $16 \times 4 = 64$  个，故优化前后均需 64 块。

input 缓存

output 缓存

列0 列1 列2 列3 列0 列1 列2 列3

|    |      |      |      |      |      |      |      |      |
|----|------|------|------|------|------|------|------|------|
| 行0 | miss |
| 行1 | hit  | hit  | miss | miss | hit  | hit  | miss | miss |
| 行2 | hit  | hit  | miss | miss | hit  | hit  | miss | miss |
| 行3 | hit  | hit  | miss | miss | hit  | hit  | miss | miss |