

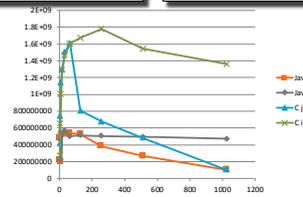
Arquitetura de Computadores

MIEI – 2017/18
DI-FCT/UNL
Aula 19

Desempenho (somadas/seg)

```
/* ij */
sum = 0.0;
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    sum += a[i][j];
  }
}

/* ji */
sum = 0.0;
for (j=0; j<n; j++) {
  for (i=0; i<n; i++) {
    sum += a[i][j];
  }
}
```

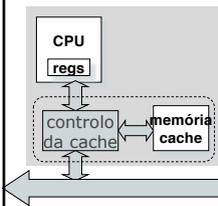


AC - 2017/18

2

Cache de memória

- Funcionamento:



CPU referencia memória:

se está em cache
accede à cache
senão
accede à memória
e atualiza a cache

A localidade espacial leva a trazer logo um bloco da memória

AC - 2017/18

3

Caching na hierarquia de memória



Nível L:
Uma memória mais rápida (mais cara, mais pequena) guarda uma cópia de alguns blocos do nível L+1

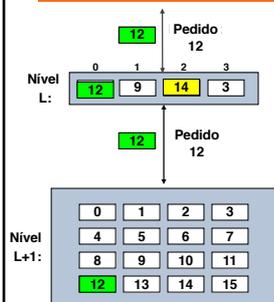
Os dados são copiados entre níveis em blocos sempre do mesmo tamanho

Nível L+1:
Uma memória maior (mais barata, mais lenta)

AC - 2017/18

4

Conceitos gerais da cache



- Um programa requer um byte de determinado bloco.
- Cache hit**
 - O programa encontra o byte na cache a nível L. E.g., bloco 14.
- Cache miss**
 - O byte não está no nível L, então deve obtê-lo no nível L+1. E.g., bloco 12.
 - Se o nível L está cheio, então é preciso escolher uma vítima:
 - Qual o bloco a ser escolhido como vítima?

AC - 2017/18

5

Conceitos gerais da cache: hit/miss

- Cache hit**
 - O programa encontra o dado pretendido na cache (nível L). Não é preciso fazer mais nada!
- Cache miss**
 - Não encontra o que pretende no nível L, então a cache de nível L deve obtê-lo no nível L+1 (e assim sucessivamente) e guardar uma cópia
 - Se o nível L está cheio, então é preciso arranjar espaço, **escolhendo uma vítima**:
 - Qual o bloco a ser escolhido como vítima (ao acaso? o menos usado? outro?)
 - Se o bloco vítima está *limpo*, isto é, **não foi alterado** desde que veio para o nível L, ou já foi atualizado em L+1, carregar o novo bloco por cima
 - Se o bloco vítima está *sujo*, isto é, **foi alterado**, é preciso escrevê-lo no nível L+1 antes de o substituir

AC - 2017/18

6

Conceitos gerais duma cache

- Taxa de sucesso (*hit ratio*)
 - h = percentagem de vezes que o dado pretendido está na cache:

$$h = \text{núm. cache hit} / \text{núm. total acessos}$$
- Tempo de acesso médio:
 - T_L é o tempo de acesso no nível L
 - T_{L+1} é o tempo de acesso no nível L+1
 - Tempo de acesso médio: T_a

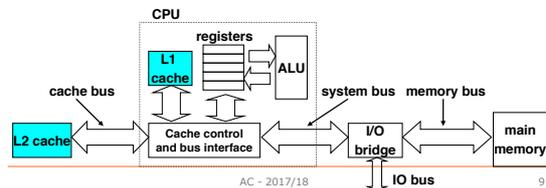
$$T_a = h * T_L + (1 - h) * T_{L+1}$$

Exemplo:

- $T_{\text{cache}} = 2 \text{ ns}$
- $T_{\text{mem}} = 8 \text{ ns}$
- *Hit ratio* = 80%
- $T_a = h * T_L + (1 - h) * T_{L+1}$
 - $T_a = 0,8 * 2 + 0,2 * 8 = 3,2 \text{ ns}$

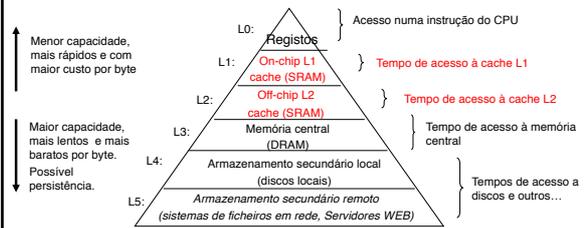
Cache – Exemplo com dois níveis

- Cache são memórias rápidas (SRAM) geridas automaticamente pelo *hardware*.
- CPU procura em L1, depois em L2, finalmente na memória central
 - Podem introduzir algum *overhead*, mas desprezável
- Exemplo de arquitetura com L2 externa ao chip do CPU:



Hierarquia de níveis de memória

- O tempo de acesso a memória não é sempre o mesmo
- O mapeamento dos dados em memória e o padrão de acesso são determinantes



Endereços e blocos de memória

Exemplo:
endereços de 8 bits
blocos de 4bytes

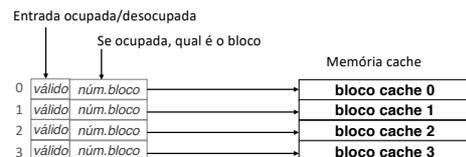
0 = 00000000	
1 = 00000001	bloco 0
2 = 00000010	
3 = 00000011	
4 = 00000100	
5 = 00000101	bloco 1
6 = 00000110	
7 = 00000111	
8 = 00010000	
9 = 00010001	bloco 2
10 = 00010010	
11 = 00010011	
12 = 00010100	
13 = 00010101	bloco 3
14 = 00010110	
15 = 00010111	
16 = 00100000	
17 = 00100001	bloco 4
18 = 00100010	
19 = 00100011	
20 = 00101000	

os 2 bits menos significativos indicam bytes dentro do mesmo bloco
os restantes correspondem ao número do bloco

■ ■ ■ ■ ■ ■ ■ ■	
núm. bloco	deslocamento no bloco

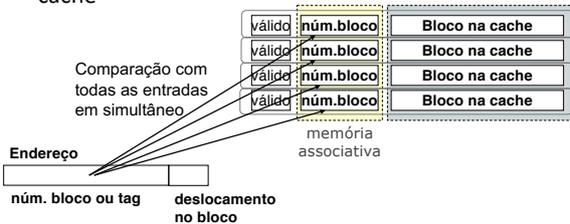
Gestão da Cache

- É necessário de forma eficiente:
 - Saber que blocos da cache estão ocupados e com que blocos de memória
 - Exemplo de cache com capacidade para 4 blocos de memória:



Cache associativa pura

- Usa o número do bloco como **chave** (ou *tag*) na memória associativa para procurar o bloco na cache



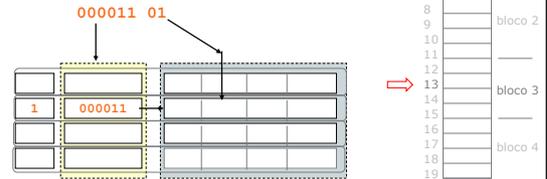
AC - 2017/18

13

Endereços e a cache

Exemplo: cache com 16bytes, 4 linhas de 4 bytes
Ler endereço 13

Bloco memória = $13/4 = 3$
 Byte nesse bloco = $13\%4 = 1$



AC - 2017/18

14

Caches associativas puras

- As Caches Associativas são muito eficientes:
 - qualquer bloco de memória pode estar em qualquer linha da cache
 - a pesquisa do bloco faz-se em paralelo, com tantos comparadores quantos os blocos que cabem na cache
- Mas são complexas e muito caras:
 - cada linha da cache deve guardar todos os bits do número do bloco nessa posição
 - cada comparador é de tantos bits quantos os bits no número do bloco
 - muitos comparadores, um para cada linha da cache

AC - 2017/18

15

Simplificando

- Pré-definindo uma linha da cache para cada bloco de memória:
 - Linha = $n^\circ \text{bloco} \% n^\circ \text{de linhas da cache}$
 - Exemplo: 4 linhas --> 2bits

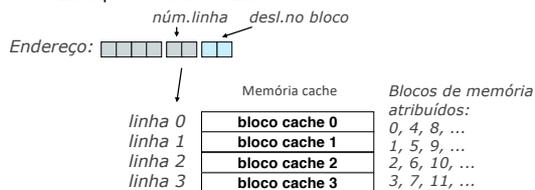
0 = 00000000	
1 = 00000001	
2 = 00000010	bloco 0
3 = 00000011	
4 = 00000100	
5 = 00000101	
6 = 00000110	bloco 1
7 = 00000111	
8 = 00010000	
9 = 00010001	
10 = 00010010	bloco 2
11 = 00010011	
12 = 00010100	
13 = 00010101	bloco 3
14 = 00010110	
15 = 00010111	
16 = 00011000	
17 = 00011001	bloco 4
18 = 00011010	
19 = 00011011	
20 = 00011100	bloco 5

AC - 2017/18

16

Simplificando

- Pré-definindo uma linha da cache para cada bloco de memória:
 - Linha = $n^\circ \text{bloco} \% n^\circ \text{de linhas da cache}$
 - Exemplo: 4 linhas --> 2bits

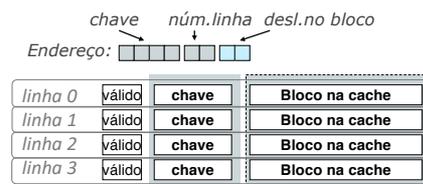


AC - 2017/18

17

Cache de mapa direto

- É necessário manter para cada linha uma chave (ou *tag*) que identifique qual o bloco na cache de entre os vários possíveis:
 - o que distingue os vários blocos são os bits mais significativos restantes do endereço
- Exemplo: end. de 8bits, cache de 4 linhas e 4 bytes p/bloco:



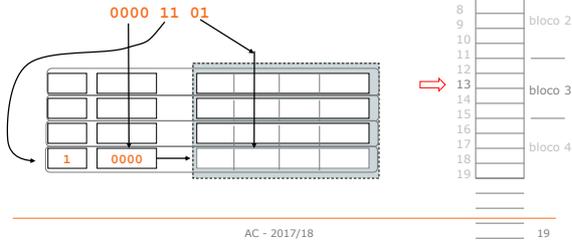
AC - 2017/18

18

Endereços e a cache

Exemplo: cache com 16bytes, 4 linhas de 4 bytes
Ler endereço 13

Bloco memória = $13/4 = 3$
 Byte nesse bloco = $13\%4 = 1$



Características da Cache de mapa direto

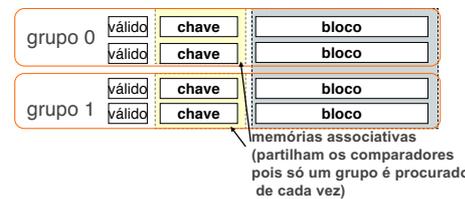
- As Caches de mapa direto são mais baratas do que as Associativas:
 - O mapa da Cache tem menos bits;
 - Não necessita de memória associativa e só precisa de 1 comparador (de tantos bits quantos os na chave dos blocos)
- São eficientes na pesquisa em Cache:
 - usam bits do endereço como índice na Cache e restantes como chave.
- Levam a colisões, mesmo que existam linhas livres!

Cache associativa por grupos (set associative)

- Solução híbrida de compromisso
- Procura-se responder às desvantagens das duas técnicas anteriores:
 - evitar as colisões no mapa direto
 - ter menos comparadores que a associativa pura
- Abordagem:
 - cada "linha" de mapa direto dá lugar a um grupo (conjunto) de blocos
 - cada bloco pode ficar em qualquer linha dentro do grupo
 - a busca no grupo é efectuada procurando a chave do bloco usando memória associativa

Cache associativa por grupos

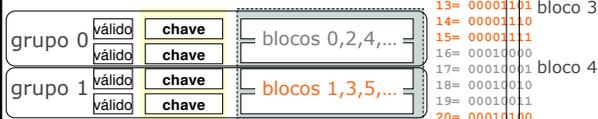
- Os blocos de memória são mapeados diretamente num grupo (ou conjunto de linhas)
 - o núm. do grupo é dado pelo: $\text{núm.bloco} \% \text{núm.grupos}$
- Dentro do grupo pode ficar em qualquer posição, como se o grupo fosse uma "mini cache associativa pura"



Cache associativa por grupos (ou conjuntos)

- Exemplo, end. de 8 bits:**
 - cache com 2 grupos
 - 4 blocos de 4 bytes, 2 por grupo
 - núm. do grupo é dado por: $\text{núm.bloco} \% \text{núm.grupos}$ (como no mapa direto para obter a linha)
 - ou seja:

chave	núm.grupo	desl.no bloco
Endereço: [][][][][][][][]		



Quando falha a cache (miss)

- Existe misses quando as localidades não se verificam:
 - por acesso a dados não contíguos em memória
 - ou estruturas que não cabem na cache
 - muda o conjunto de instruções em execução (*working-set*): devido a jumps, calls, ...
 - ou este conjunto não cabe todo na cache
 - o SO troca de programa em execução
 - muda de código e de dados...