Arquitetura de Computadores

MiEI – 2016/17 DI-FCT/UNL Aula 18

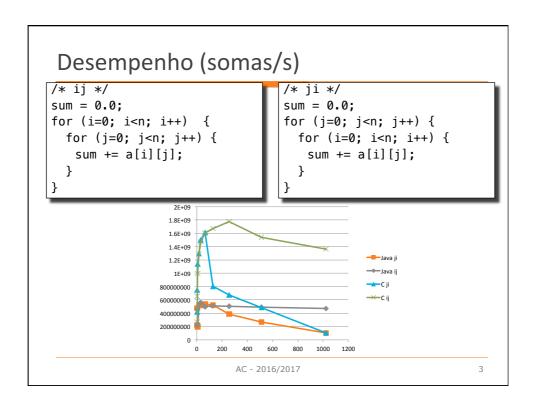
AC - 2016/2017

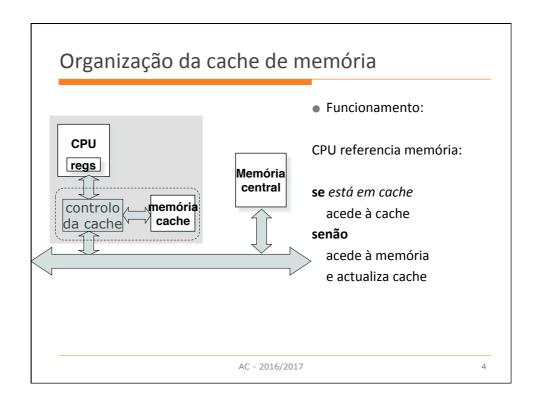
1

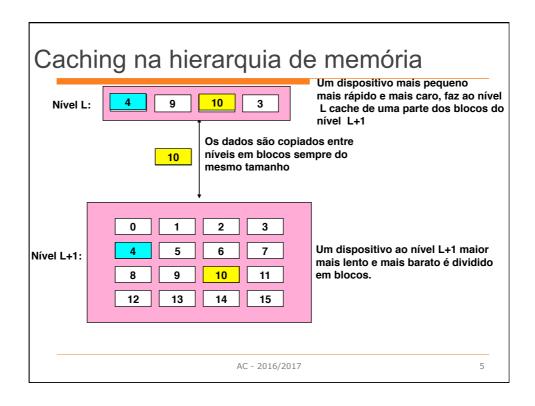
Ideia da Cache

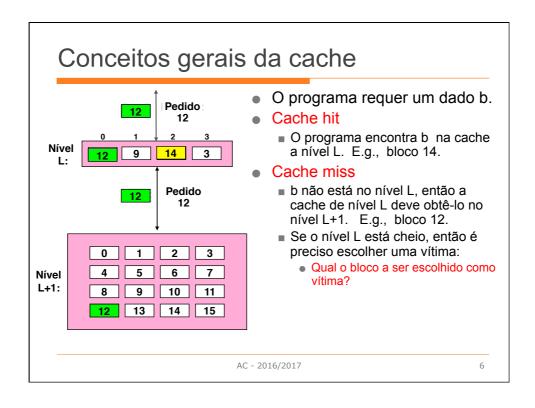
- Vantagens esperadas:
 - ■A localidade leva a que a maior parte dos acessos (leitura/escrita) sejam aos dados já no nível L em vez de ter de ao nível L+1.
 - ■O armazenamento no nível L+1 pode ser mais lento, e assim mais barato e muito maior.
 - ■Efeito global: Um grande conjunto de memória que custa o preço da que se encontra no nível mais baixo, mas que permite o acesso aos dados a um ritmo semelhante à do nível mais elevado.

AC - 2016/2017









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 enconta o que pretende no nível L, então a cache de nível L deve obtê-lo no nível L+1 (e assim sucessivamente).
- Se o nível L está cheio, então é preciso arranjar espaço, escolhendo uma vítima:
 - Qual o bloco a ser escolhido como vítima (ao acaso? LRU least recently used?)
 - Se o bloco vítima está limpo, isto é, não foi alterado desde que veio para o nível L, posso 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 - 2016/2017

7

Conceitos gerais duma cache

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

h = núm. cache hit / núm. total acessos

- Tempo de acesso médio:
 - lacksquare T_L é o tempo de acesso no nível L
 - T_{I+1} é o tempo de acesso no nível L+1
 - Tempo de acesso médio: Ta

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

AC - 2016/2017

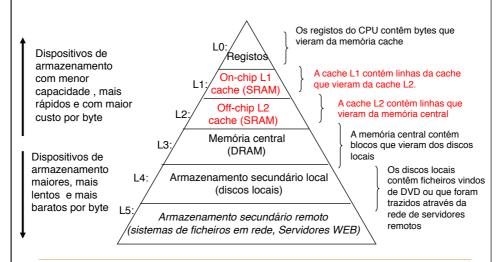
• Exemplo:

- T cache = 2 ns
- $\blacksquare T mem = 8 ns$
- **■***Hit ratio* = 80%
- Ta = $h * T_L + (1 h) * T_{L+1}$ • Ta = 0,8 x 2 + 0,2 x 8 = **3,2** ns

AC - 2016/2017

9

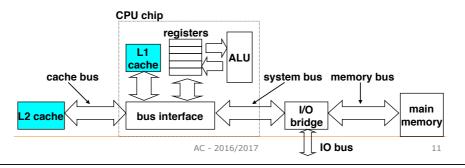
Hierarquia de níveis de memória



AC - 2016/2017

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:



Política de escritas com cache

- Quando o CPU altera o conteúdo duma posição de memória, faz sentido que a alteração fique em cache
 - ■a ideia é que o CPU pode vir ler esse valor em breve (localidade temporal)
- O que acontece na memória de nível inferior, quando um bloco na cache é alterado, depende da política de escrita:
 - ■Write-through
 - ■Write-back (ou delayed)

AC - 2016/2017

Write-through

- A memória cache e a memória central são actualizadas em cada escrita
- Assegura que a cache e a memória central estão sempre coerentes
- Cada escrita demora o tempo de escrita na memória central
- Tolerável, porque normalmente há muito mais acessos em leitura do que em escrita

AC - 2016/2017

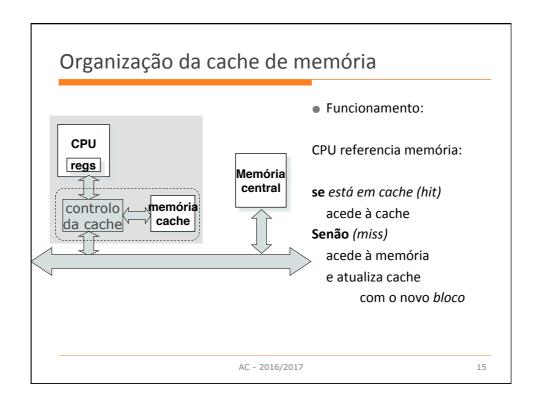
13

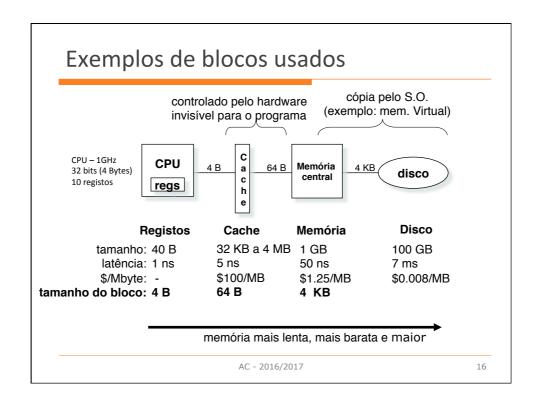
Write-back

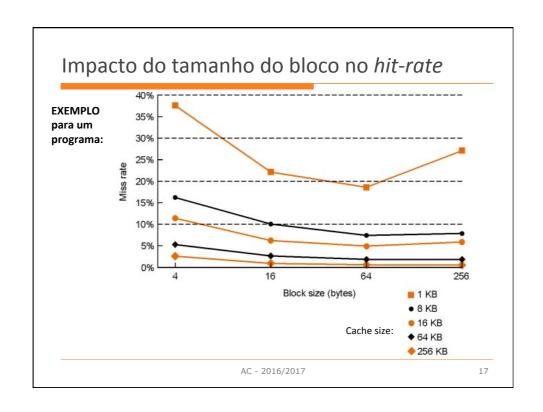
- Quando há uma escrita só é atualizada a cache; mas marca-se esse bloco como alterado (dirty)
- Só quando um bloco é removido da cache é que se atualiza o bloco na memória de nível inferior
- Isto é mais rápido do que o write-through
- Diminui o tráfego para a memória de nível inferior
- O inconveniente é que a cache e a memória inferior nem sempre estão coerentes

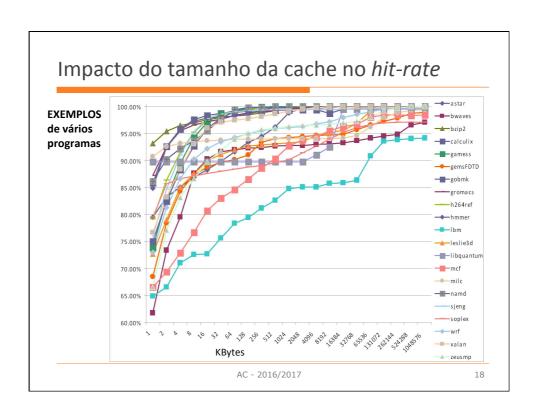
■isto pode provocar problemas ...

AC - 2016/2017









Organização da cache

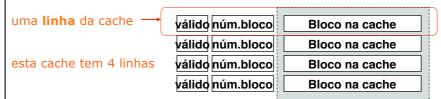
- Para além de guardar o conteúdo da memória de nível inferior, também necessita de guardar:
 - •Que zonas da cache estão realmente em uso (de inicio a cache está vazia!)
 - •Que endereços/blocos da memória de nível inferior estão na cache
 - ■Se usar write-back, que zonas foram escritas
 - ■Informação que permita escolher os blocos a substituir
- A consulta e atualização desta informação tem de ser eficiente!!!

AC - 2016/2017 19

Endereços e blocos de memória 0= 00000000 Exemplo: 1 = 000000 01 bloco 02= 00000010 endereços de 8 bits 3= 00000011 blocos de 4bytes 4= 00000100 5= 000001 01 bloco 1 6= 00000110 dado um endereço E: 7= 00000111 E/4 = número do bloco 8= 00001000 bloco 2 E%4 = byte dentro do bloco 9= 00001001 ou seja: 10= 00001010 11= 00001011 4 bytes por bloco = $2^2 \rightarrow 2$ bits 12= 00001100 13= 00001101 bloco 3 14= 00001110 os 2 bits menos significativos indicam bytes dentro do mesmo bloco 15= 00001111 16= 00010000 os restantes correspondem ao número do bloco 17= 00010001 bloco 4 18= 00010010 núm. bloco deslocamento no bloco 19= 00010011 20= 00010100 AC - 2016/2017 20

Exemplo (write-through)

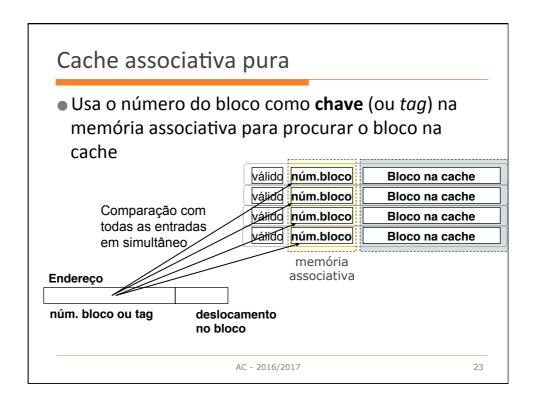
 A cache tem para cada bloco, se está ocupado (válido) ou não e o respectivo número de bloco, quando ocupado

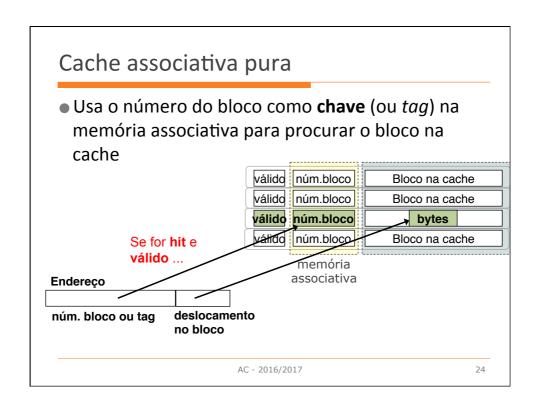


- Para cada acesso do CPU à memória, parte-se o endereço de acordo com o tamanho do bloco e procura-se o seu número na cache
 - ■só faz sentido se esta procura for rápida!
- Usa-se os bits de deslocamento no bloco para aceder ao byte pretendido

AC - 2016/2017 21

Memória associativa (interrogável pelo conteúdo) Exemplo com 4 entradas valor a procurar 110011 entrada 0 010101 entrada 1 110011 comparadores entrada 2 010000 entrada 3 111011 Compara em paralelo com todas as entradas hit na entrada 1 AC - 2016/2017 22

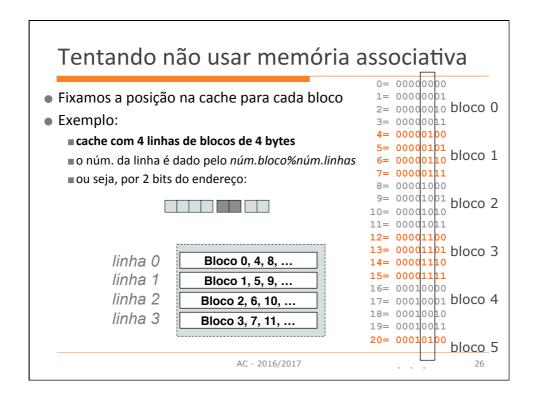




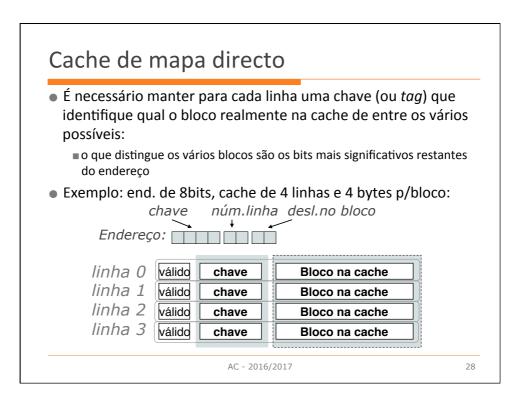
Caches associativas puras

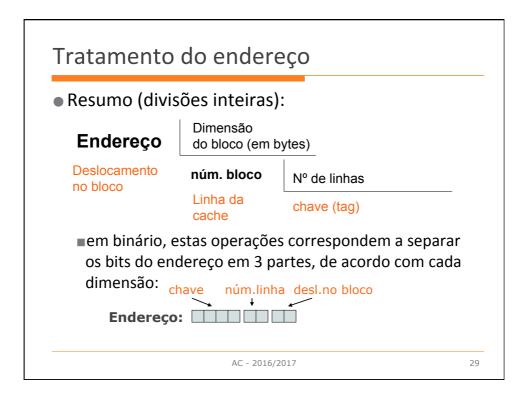
- As Caches Associativas são muito eficientes:
 - ■qualquer bloco pode estar em qualquer posição
 - a pesquisa do bloco faz-se em paralelo, com tantos comparadores quantos os blocos que cabem na cache
- Mas são complexas e muito caras:
 - ■cada linha da cache deve guardar todos os bits do número do bloco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

AC - 2016/2017



Porque se usam os bits do meio para				
escolher a linha?				
Cache com 4 linhas	Utilizar os bits Utilizar os bits Mais significativos do meio			
00	<u>00</u> 00		00 <u>00</u>	
01	<u>00</u> 01		00 <u>01</u>	
10	<u>00</u> 10		00 <u>10</u>	
11	<u>00</u> 11		00 <u>11</u>	
Se usasse os bits mais signif.:	<u>01</u> 00		01 <u>00</u>	
■ Blocos de memória adjacentes	<u>01</u> 01		01 <u>01</u>	
corresponderiam à mesma linha na	<u>01</u> 10		01 <u>10</u>	
cache	<u>01</u> 11		01 <u>11</u>	
	<u>10</u> 00		10 <u>00</u>	
■ Mau para a localidade espacial!	<u>10</u> 01		10 <u>01</u>	
Usando os bits do meio:	<u>10</u> 10		10 <u>10</u>	
■ Blocos consecutivas de memória	<u>10</u> 11		10 <u>11</u>	
correspondem a linhas da cache	<u>11</u> 00		11 <u>00</u>	
distintas	<u>11</u> 01		11 <u>01</u>	
	<u>11</u> 10		11 <u>10</u>	
■ Melhor para a localidade espacial!	<u>11</u> 11		11 <u>11</u>	
AC - 2016/2017 27				





Vantagens da Cache de mapa directo

- 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:
 - ■porque usam os bits do endereço diretamente como índice no mapa da Cache e como chave, só precisando de comparar essa chave.

AC - 2016/2017

Desvantagems da Cache de mapa directo

- Levam a colisões: todos os blocos com a mesma linha atribuída colidem nessa posição
 - mesmo que existam outras linhas na cache livres!
- A probabilidade de colisões diminui à medida que se aumenta a Capacidade da Cache
 - ■mas o problema mantém-se ...

AC - 2016/2017