

# Arquitetura de Computadores

MIEI – 2018/19  
DI-FCT/UNL  
Aula 20

AC - 2018/19

1

## Princípio da Localidade

- Os programas têm padrões de acesso a memória:
  - Localidade temporal**: itens referenciados recentemente tendem a voltar a sê-lo.
  - Localidade espacial**: itens com endereços próximos tendem a ser referenciados em conjunto.

Exemplo:

**Dados**

Elementos do array são referenciados em sequência:  
sum referenciado em cada iteração:

**Instruções**

Instruções referenciadas em sequência:

Ciclo :

```
sum = 0;  
for (i = 0; i < n; i++)  
    sum += a[i];  
return sum;
```

Localidade espacial

Localidade temporal

AC - 2018/19

2

## Hierarquia de níveis de memória

- Poderemos ter mais um nível para que o programa que está a executar seja mais rápido?

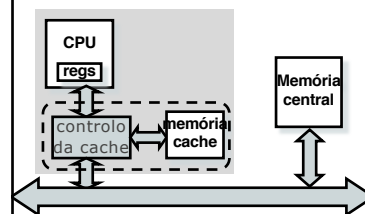


AC - 2018/19

3

## Cache de memória

- Funcionamento:



CPU referencia memória:

se está em cache

accede à cache

senão

accede à memória

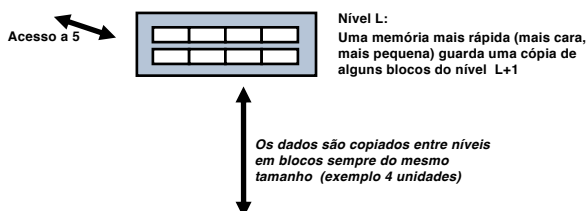
e atualiza a cache

A localidade espacial leva a trazer logo um bloco da memória

AC - 2018/19

4

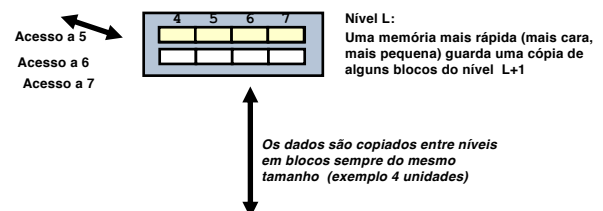
## Caching na hierarquia de memória



AC - 2018/19

5

## Caching na hierarquia de memória



AC - 2018/19

6

## Conceitos gerais de funcionamento

- Para cada acesso à memória de determinado bloco:
- **Cache hit**
  - Este está na cache → tempo de acesso rápido (nível L)
  - Exemplo anterior: acesso aos endereços 6 e 7
- **Cache miss**
  - Este não está na cache. Deve obtê-lo no nível inferior → tempo de acesso lento (nível L+1)
  - Exemplo anterior: primeiro acesso, ao endereço 5
  - Se a cache está cheia, então é preciso arranjar espaço:
    - Qual o bloco a eliminar?

AC - 2018/19

7

## Tratamento do miss

- Cache miss
  - Não encontra o que pretende no nível L, então vamos ao nível L+1 (e assim sucessivamente)
  - Guardar uma cópia do bloco no nível L
  - Mas se o nível L está cheio, então é preciso arranjar espaço:
    - Que bloco escolher como vítima? Ao acaso? O menos usado? ...?
    - Se o bloco vítima não foi alterado, carregar o novo bloco por cima
    - Se o bloco vítima está alterado (foi escrito), é preciso garantir que no nível L+1 também está alterado antes de o substituir

AC - 2018/19

8

## Tempo médio de acesso

- 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 médio de acesso:
  - $T_L$  é o tempo de acesso no nível L
  - $T_{L+1}$  é o tempo de acesso no nível L+1
  - Tempo médio de acesso,  $T_a$
  - $T_a = h * T_L + (1 - h) * T_{L+1}$

AC - 2018/19

9

## Exemplo:

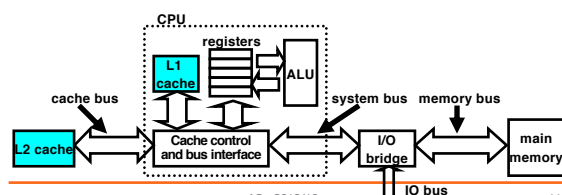
- $T_{\text{cache}} = 2 \text{ ns}$
- $T_{\text{mem}} = 8 \text{ ns}$
- $\text{Hit ratio} = 80\%$  numa determinada execução
- $T_a = h * T_L + (1 - h) * T_{L+1}$
- $T_a = 0,8 * 2 + 0,2 * 8 = 3,2 \text{ ns}$
- Nesta execução a percepção do tempo de acesso à memória é de 3,2 ns

AC - 2018/19

10

## Cache – Exemplo com dois níveis

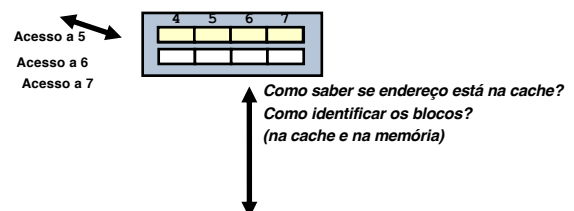
- Cache são memórias rápidas (SRAM) geridas automaticamente pelo *hardware*.
- CPU procura em L1, depois em L2, finalmente na memória central
  - Podem introduzir algum *overhead*, mas desprezável
- Exemplo de arquitetura com L2 externa ao chip do CPU:



AC - 2018/19

11

## Caching na hierarquia de memória

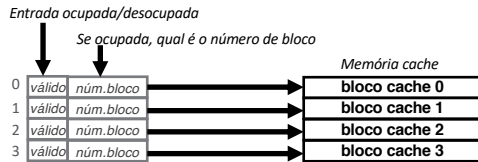


AC - 2018/19

13

## 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:



AC - 2018/19

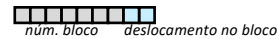
14

## Endereços e blocos de memória

Exemplo:  
endereços de 8 bits  
blocos de 4 bytes

dado um endereço E:  
 $E/4$  = número do bloco  
 $E\%4$  = byte dentro do bloco  
 ou seja:  
 4 bytes por bloco =  $2^2 \rightarrow 2$  bits

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



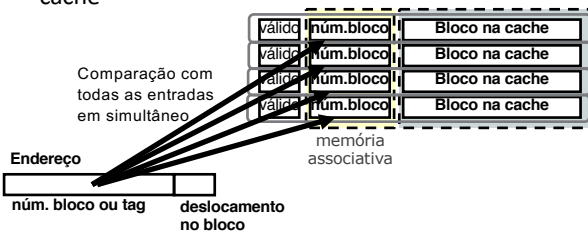
0=	00000000	
1=	00000001	bloco 0
2=	00000010	
3=	00000011	
4=	00000100	bloco 1
5=	00000101	
6=	00000110	
7=	00000111	
8=	00001000	bloco 2
9=	00001001	
10=	00001010	
11=	00001011	
12=	00001100	bloco 3
13=	00001101	
14=	00001110	
15=	00001111	
16=	00010000	bloco 4
17=	00010001	
18=	00010010	
19=	00010011	
20=	00010100	

AC - 2018/19

15

## Cache associativa pura

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



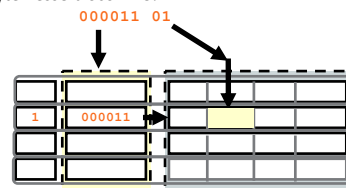
AC - 2018/19

16

## Endereços e a cache

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

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



0																			
1																			
2																			
3																			
4																			
5																			
6																			
7																			
8																			
9																			
10																			
11																			
12																			
13																			
14																			
15																			
16																			
17																			
18																			
19																			

AC - 2018/19

17

## 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 - 2018/19

18

## 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  $\rightarrow$  2bits

0=	00010000	
1=	00010001	
2=	00010010	bloco 0
3=	00010011	
4=	00010100	
5=	00010101	bloco 1
6=	00010110	
7=	00010111	
8=	00011000	
9=	00011001	bloco 2
10=	00011010	
11=	00011011	
12=	00011100	
13=	00011101	bloco 3
14=	00011110	
15=	00011111	
16=	00010000	
17=	00010001	bloco 4
18=	00010010	
19=	00010011	
20=	00010100	bloco 5

AC - 2018/19

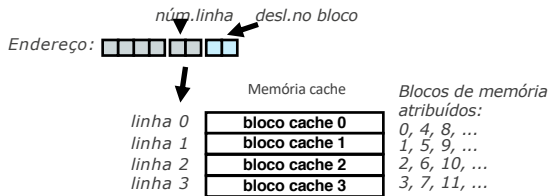
19

## Simplificando

- Pré-definindo uma linha da cache para cada bloco de memória:

■  $\text{Linha} = \text{n}^\circ \text{bloco} \% \text{n}^\circ \text{de linhas da cache}$

■ Exemplo: 4 linhas --> 2bits

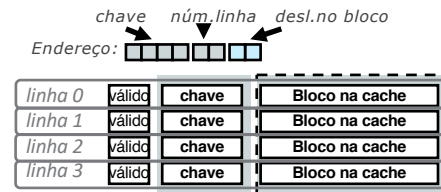


AC - 2018/19

20

## 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 - 2018/19

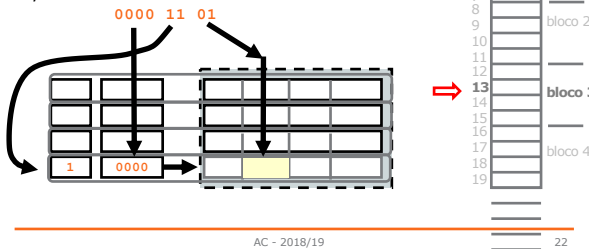
21

## 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 - 2018/19

22

## 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!

AC - 2018/19

23

## Cache associativa por grupos (set associative)

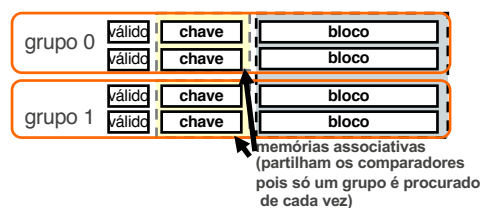
- Solução híbrida de compromisso
- Procura-se responder às desvantagens das duas técnicas anteriores:
  - evitar as colisões no mapa direto
  - ter menos comparadores que a associativa pura
- Abordagem:
  - cada "linha" de mapa direto dá lugar a um grupo (conjunto) de blocos
  - cada bloco pode ficar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

AC - 2018/19

24

## 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"



AC - 2018/19

25

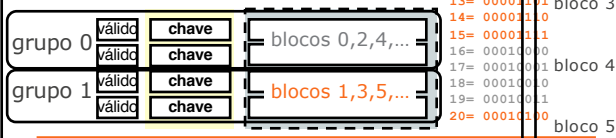
## Cache associativa por grupos (ou conjuntos)

- 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:



AC - 2018/19

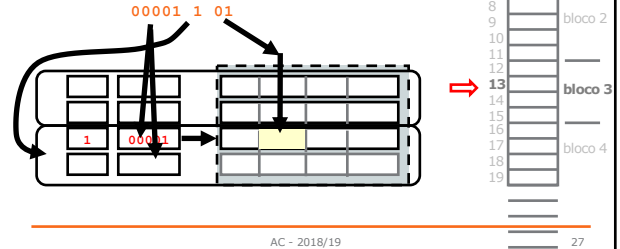
26

## 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 - 2018/19

27

## 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...

AC - 2018/19

28

## Tratamento dos misses

- Se **Cache miss**: é preciso arranjar espaço para o novo bloco. Se associativa por grupos o conflito tem de ser resolvido apenas no grupo.
- Se linha disponível (na cache ou no grupo), marca como válida, atualiza a chave (tag) e o conteúdo do bloco
- Se cache cheia, qual o bloco a ser eliminado da cache ?
  - Um qualquer . . . ou o bloco que no futuro não vai ser necessário. Mas qual é esse?
- Se usando *write-back*, verifica bit *dirty* e atualiza memória, antes de usar a linha vítima

AC - 2018/19

29