

## Organização de Computadores

### Aula 16

### Memória cache primeira parte

INF01113 - Organização de Computadores

## Memória cache primeira parte

1. Tendências tecnológicas
2. Hierarquia de memória
3. Princípio de localidade
4. Impacto no desempenho
5. Organizações de memória cache

INF01113 - Organização de Computadores

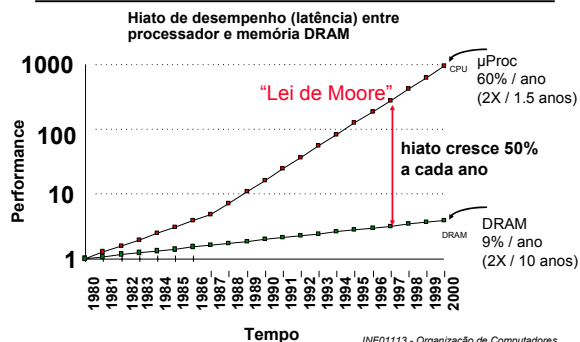
### 1. Tendências tecnológicas

	Capacidade	Velocidade (latência)
Lógica:	2x em 3 anos	2x em 3 anos
DRAM:	4x em 3 anos	2x em 10 anos
Disco:	4x em 3 anos	2x em 10 anos

DRAM		
Ano	Tamanho	Tempo acesso
1980	64 Kb	250 ns
1983	256 Kb	220 ns
1986	1 Mb	190 ns
1989	4 Mb	165 ns
1992	16 Mb	145 ns
1995	64 Mb	120 ns

INF01113 - Organização de Computadores

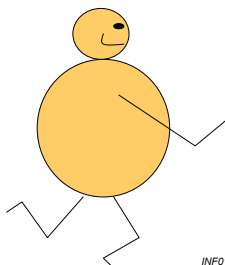
### Tendências tecnológicas



INF01113 - Organização de Computadores

### O que queremos?

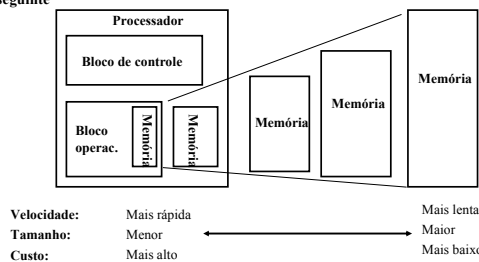
- Um lutador de sumo que corre e ganha os 100m rasos!



INF01113 - Organização de Computadores

### 2. Hierarquia de memória

- objetivo: oferecer ilusão de máximo tamanho de memória, com mínimo custo e máxima velocidade
- cada nível contém cópia de parte da informação armazenada no nível superior seguinte



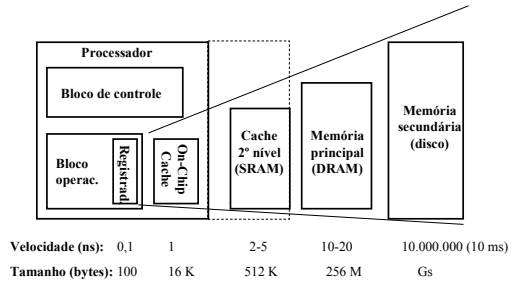
INF01113 - Organização de Computadores

## Tecnologias na hierarquia de memória

- Acesso randômico
  - tempo de acesso é o mesmo para todas as posições
  - **DRAM**: Dynamic Random Access Memory
    - alta densidade, baixa potência, barata, lenta
    - dinâmica: precisa de um “refresh” regular
  - **SRAM**: Static Random Access Memory
    - baixa densidade, alta potência, cara, rápida
    - estática: conteúdo dura “para sempre”(enquanto houver alimentação)
- Acesso “não-tão-randômico”
  - tempo de acesso varia de posição para posição e de tempos em tempos
  - exemplos: disco, CD-ROM
- Acesso sequencial
  - tempo de acesso varia linearmente com a posição (p.ex. fita)

INF01113 - Organização de Computadores

## Hierarquia de memória



INF01113 - Organização de Computadores

## Hit e miss

- **Hit**: dado aparece em algum bloco no nível superior (junto ao processador)
  - **Hit Ratio**: a fração de acessos à memória resolvidos no nível superior
  - **Hit Time**: tempo de acesso ao nível superior, que consiste de tempo de acesso à memória RAM + tempo para determinar hit/miss
- **Miss**: dado precisa ser buscado de um bloco no nível inferior
  - **Miss Ratio** = 1 – (Hit Ratio)
  - **Miss Penalty**: tempo gasto para substituir um bloco no nível superior + tempo para fornecer o bloco ao processador
- Hit Time << Miss Penalty

INF01113 - Organização de Computadores

## Hierarquia de memória

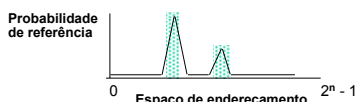
Como a hierarquia é gerenciada?

- Registradores <=> memória
  - pelo compilador
- cache <=> memória principal
  - pelo hardware
- memória principal <=> disco
  - pelo hardware e pelo sistema operacional (memória virtual)
  - pelo programador (arquivos)

INF01113 - Organização de Computadores

## 3. Princípio de localidade

- Hierarquia de memória funciona devido ao princípio de localidade
  - todos os programas repetem trechos de código e acessam repetidamente dados próximos



- **localidade temporal**: posições de memória, uma vez acessadas, tendem a ser acessadas novamente no futuro próximo
- **localidade espacial**: endereços em próximos acessos tendem a ser próximos de endereços de acessos anteriores

INF01113 - Organização de Computadores

## Exemplo:

- 1,2,3,4,5,6,7,6,7,6,7,492, 493,494,6,7,6,7...

INF01113 - Organização de Computadores

## Princípio de localidade

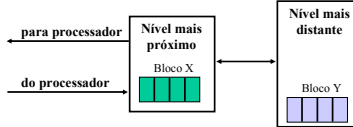
Como explorar o princípio de localidade numa hierarquia de memória?

- Localidade Temporal

=> Mantenha itens de dados mais recentemente acessados nos níveis da hierarquia mais próximos do processador

- Localidade Espacial

=> Mova blocos de palavras contíguas para os níveis da hierarquia mais próximos do processador



INF01113 - Organização de Computadores

## Localidade temporal

- usualmente encontrada em laços de instruções e acessos a pilhas de dados e variáveis
- é essencial para a eficiência da memória cache
- se uma referência é repetida N vezes durante um laço de programa, após a primeira referência a posição é sempre encontrada na cache

$T_c$  = tempo de acesso à cache

$T_m$  = tempo de acesso à memória principal

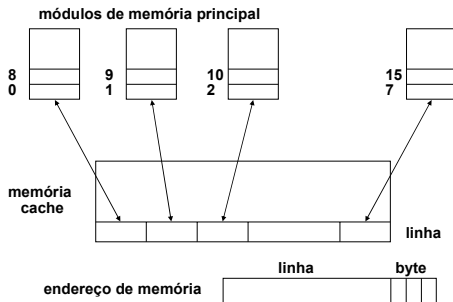
$T_{ce}$  = tempo efetivo de acesso à cache

$$T_{ce} = \frac{N T_c + T_m}{N} = T_c + \frac{T_m}{N}$$

se  $T_c = 1 \text{ ns}$ ,  $T_m = 20 \text{ ns}$ ,  $N = 10 \Rightarrow T_{ce} = 3 \text{ ns}$   
 $N = 100 \Rightarrow T_{ce} = 1,2 \text{ ns}$

INF01113 - Organização de Computadores

## Localidade espacial



INF01113 - Organização de Computadores

## Localidade espacial

- memória principal é entrelaçada
- uma linha é transferida num único acesso entre a memória principal e a cache, através de um largo barramento de dados
- casamento entre tempo de acesso da cache e da memória principal

$M$  = número de módulos da memória principal

$T_c$  = tempo de acesso à cache

$T_m$  = tempo de acesso à memória principal

Ideal:  $M T_c = T_m$

INF01113 - Organização de Computadores

## Localidade espacial

- tempo médio de acesso a um byte, na primeira referência  
 $= 2 M T_c / M = 2 T_c$
- se cada byte é referenciado N vezes na cache, então o tempo efetivo (médio) de acesso  $T_{ce}$  a cada byte é

$$T_{ce} = \frac{2 T_c + (N - 1) T_c}{N} = \frac{(N + 1) T_c}{N}$$

se  $T_c = 1 \text{ ns}$   
 $T_m = 20 \text{ ns}$   
 $N = 10$

então  $T_{ce} = 1,1 \text{ ns}$

$N = 100$   
 $T_{ce} = 1,01 \text{ ns}$

INF01113 - Organização de Computadores

## 4. Impacto no desempenho

Medindo o impacto do **hit ratio** no tempo efetivo de acesso

$T_c$  = tempo de acesso à memória cache

$T_m$  = tempo de acesso à memória principal

$T_{ce}$  = tempo efetivo de acesso à memória cache, considerando efeito dos misses

$$T_{ce} = T_c + (1 - h) T_m$$

se  $T_c = 1 \text{ ns}$ ,  $T_m = 20 \text{ ns}$

$h = 0.85 \quad 0.95 \quad 0.99 \quad 1.0$

então  $T_{ce} = 4 \text{ ns} \quad 2 \text{ ns} \quad 1.2 \text{ ns} \quad 1 \text{ ns}$

INF01113 - Organização de Computadores

## Impacto no desempenho

Tempo gasto com um *cache miss*, em número de instruções executadas

1° Alpha	340 ns / 5.0 ns = 68 clks x 2 instr.	ou	136 instruções
2° Alpha	266 ns / 3.3 ns = 80 clks x 4 instr.	ou	320 instruções
3° Alpha	180 ns / 1.7 ns = 108 clks x 6 instr.	ou	648 instruções

$$1/2 \times \text{latência} \times 3 \times \text{frequência clock} \times 3 \times \text{instruções/clock} \Rightarrow \approx 5 \times$$

INF01113 - Organização de Computadores

## Impacto no desempenho

- Supondo um processador que executa um programa com:

- CPI = 1.1
- 50% aritm/lógica, 30% load/store, 20% desvios

- Supondo que 10% das operações de acesso a dados na memória sejam *misses* e resultem numa penalidade de 50 ciclos

$$\begin{aligned} \text{CPI} &= \text{CPI ideal} + n^\circ \text{ médio de stalls por instrução} \\ &= 1.1 \text{ ciclos} + 0.30 \text{ acessos à memória / instrução} \\ &\quad \times 0.10 \text{ misses / acesso} \times 50 \text{ ciclos / miss} \\ &= 1.1 \text{ ciclos} + 1.5 \text{ ciclos} \\ &= 2.6 \end{aligned}$$

CPI ideal	1.1
Data misses	1.5
Instr.misses	0.5

- 58 % do tempo o processador está parado esperando pela memória!
- um miss ratio de 1% no fetch de instruções resultaria na adição de 0.5 ciclos ao CPI médio

INF01113 - Organização de Computadores

## 5. Organizações de memória cache

- processador gera endereço de memória e o envia à cache
- cache deve
  - verificar se tem cópia da posição de memória correspondente
  - se tem, encontrar a posição da cache onde está esta cópia
  - se não tem, trazer o conteúdo da memória principal e escolher posição da cache onde a cópia será armazenada
- mapeamento* entre endereços de memória principal e endereços de cache resolve estas 3 questões
  - deve ser executado em hardware
- estratégias de organização (mapeamento) da cache
  - mapeamento completamente associativo
  - mapeamento direto
  - mapeamento set-associativo

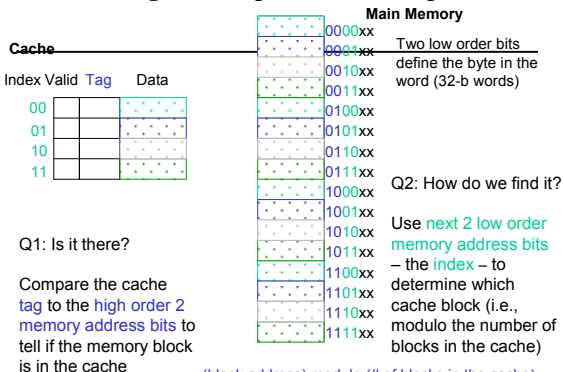
INF01113 - Organização de Computadores

## Cache

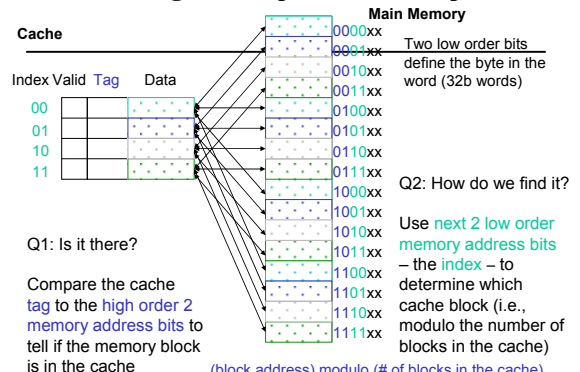
- Two questions to answer (in hardware):
  - Q1: How do we know if a data item is in the cache?
  - Q2: If it is, how do we find it?
- Direct mapped
  - For each item of data at the lower level, there is exactly one location in the cache where it might be - so lots of items at the lower level must *share* locations in the upper level
  - Address mapping:
 
$$(\text{block address}) \bmod (\# \text{ of blocks in the cache})$$
  - First consider block sizes of *one word*

INF01113 - Organização de Computadores

## Caching: A Simple First Example



## Caching: A Simple First Example



## Direct Mapped Cache

- Consider the main memory word reference string  
0 1 2 3 4 3 4 14

Start with an empty cache - all blocks initially marked as not valid

0	1	2	3																																
<table><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>									<table><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>									<table><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>									<table><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>								
4	3	4	15																																
<table><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>									<table><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>									<table><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>									<table><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>								

INF01113 - Organização de Computadores

## Direct Mapped Cache

- Consider the main memory word reference string

Start with an empty cache - all blocks initially marked as not valid  
0 1 2 3 4 3 4 15

0 miss		1 miss		2 miss		3 miss	
00	Mem(0)	00	Mem(0)	00	Mem(0)	00	Mem(0)
		00	Mem(1)	00	Mem(1)	00	Mem(1)
				00	Mem(2)	00	Mem(2)
						00	Mem(3)

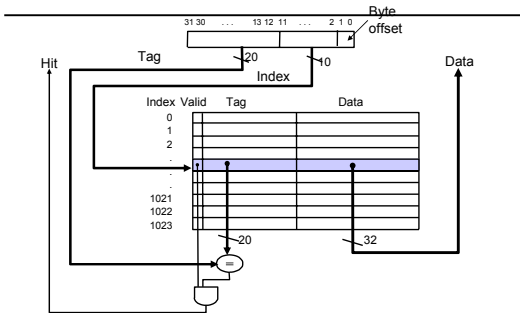
4 miss		3 hit		4 hit		15 miss	
00	Mem(0)	01	Mem(4)	01	Mem(4)	01	Mem(4)
00	Mem(1)	00	Mem(1)	00	Mem(1)	00	Mem(1)
00	Mem(2)	00	Mem(2)	00	Mem(2)	00	Mem(2)
00	Mem(3)	00	Mem(3)	00	Mem(3)	00	Mem(3)

- 8 requests, 6 misses

INF01113 - Organização de Computadores

## MIPS Direct Mapped Cache Example

- One word/block, cache size = 1K words



What kind of locality are we taking advantage of?

INF01113 - Organização de Computadores