

Organização de Computadores

Aula 16

Memória Cache primeira parte

Memória Cache

primeira parte

- 1. Tendências Tecnológicas**
- 2. Hierarquia de Memória**
- 3. Princípio de Localidade**
- 4. Impacto no Desempenho**
- 5. Organizações de Memória Cache**

1. Tendências Tecnológicas

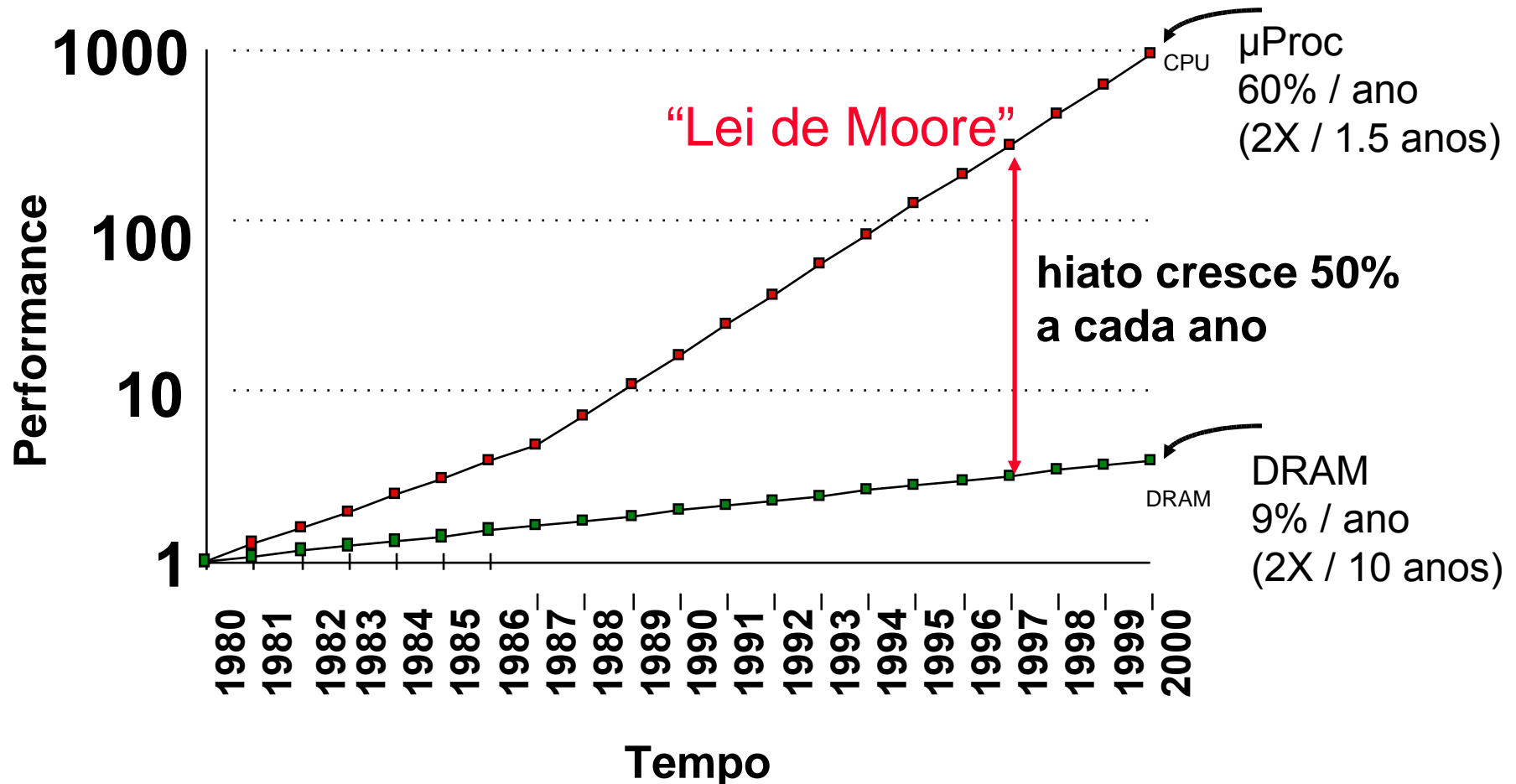
1. Tendências Tecnológicas

| | 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 | | |
|------------|----------------|---------------------|
| <u>Ano</u> | <u>Tamanho</u> | <u>Tempo acesso</u> |
| 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 |

Tendências Tecnológicas

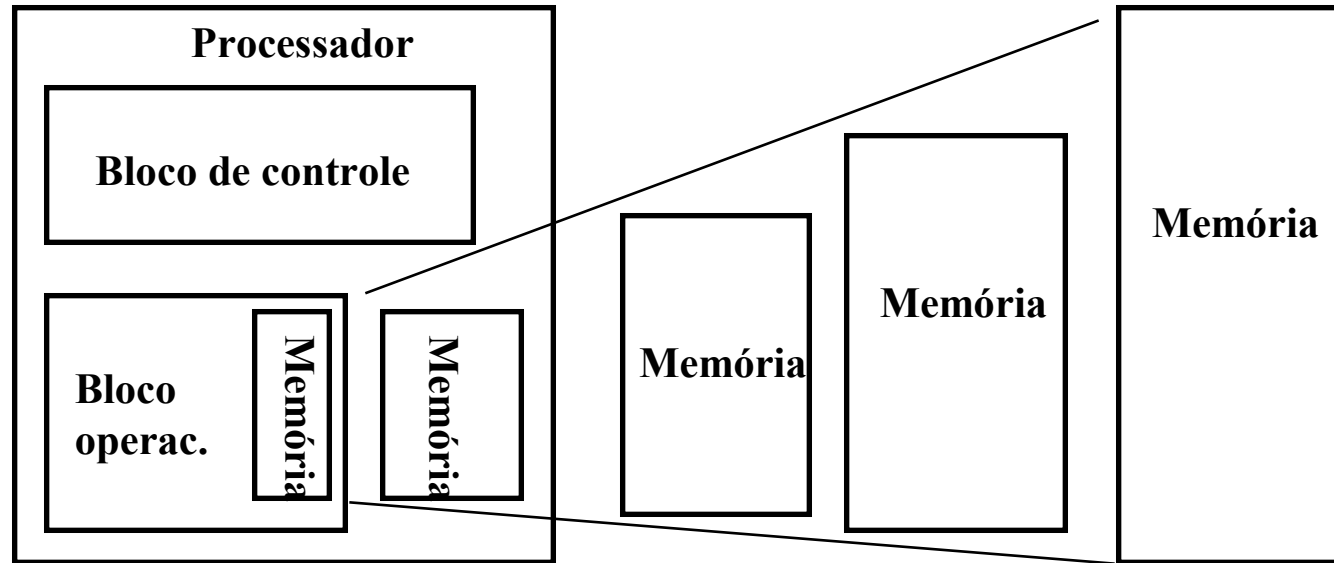
Hiato de desempenho (latência) entre
processador e memória DRAM



2. Hierarquia de Memória

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



Velocidade:

Mais rápida

Tamanho:

Menor

Custo:

Mais alto

Mais lenta

Maior

Mais baixo



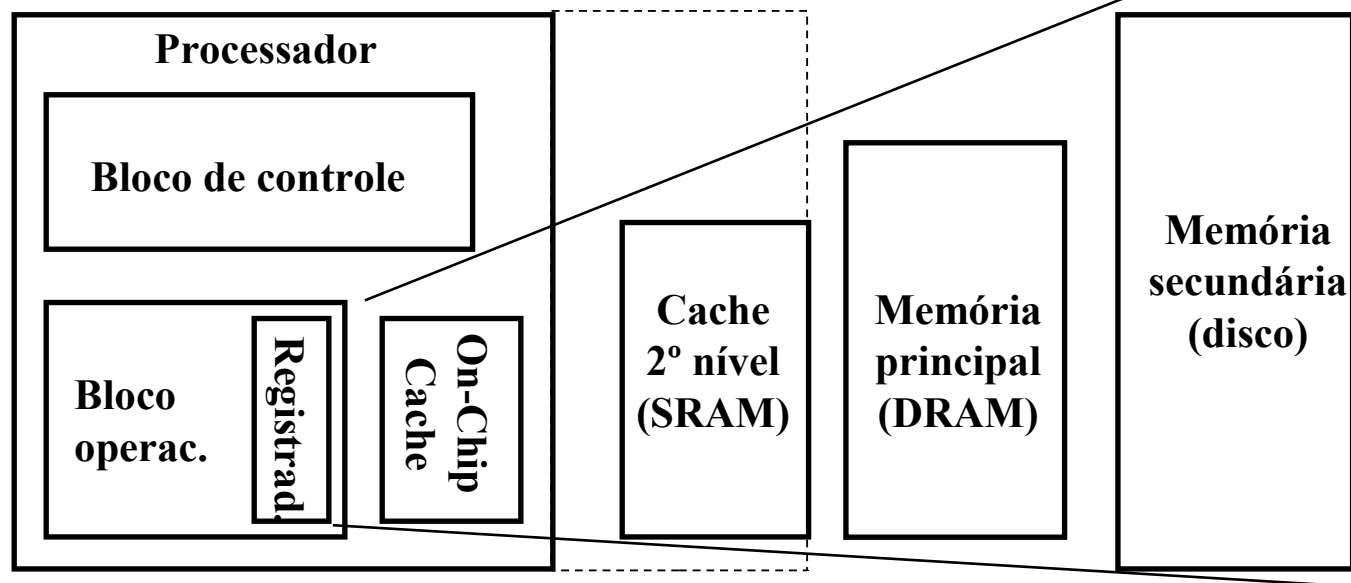
Tecnologias na Hierarquia de Memória

- **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)

Hierarquia de Memória

Perto da CPU – rápido, menor e acesso freqüente aos dados

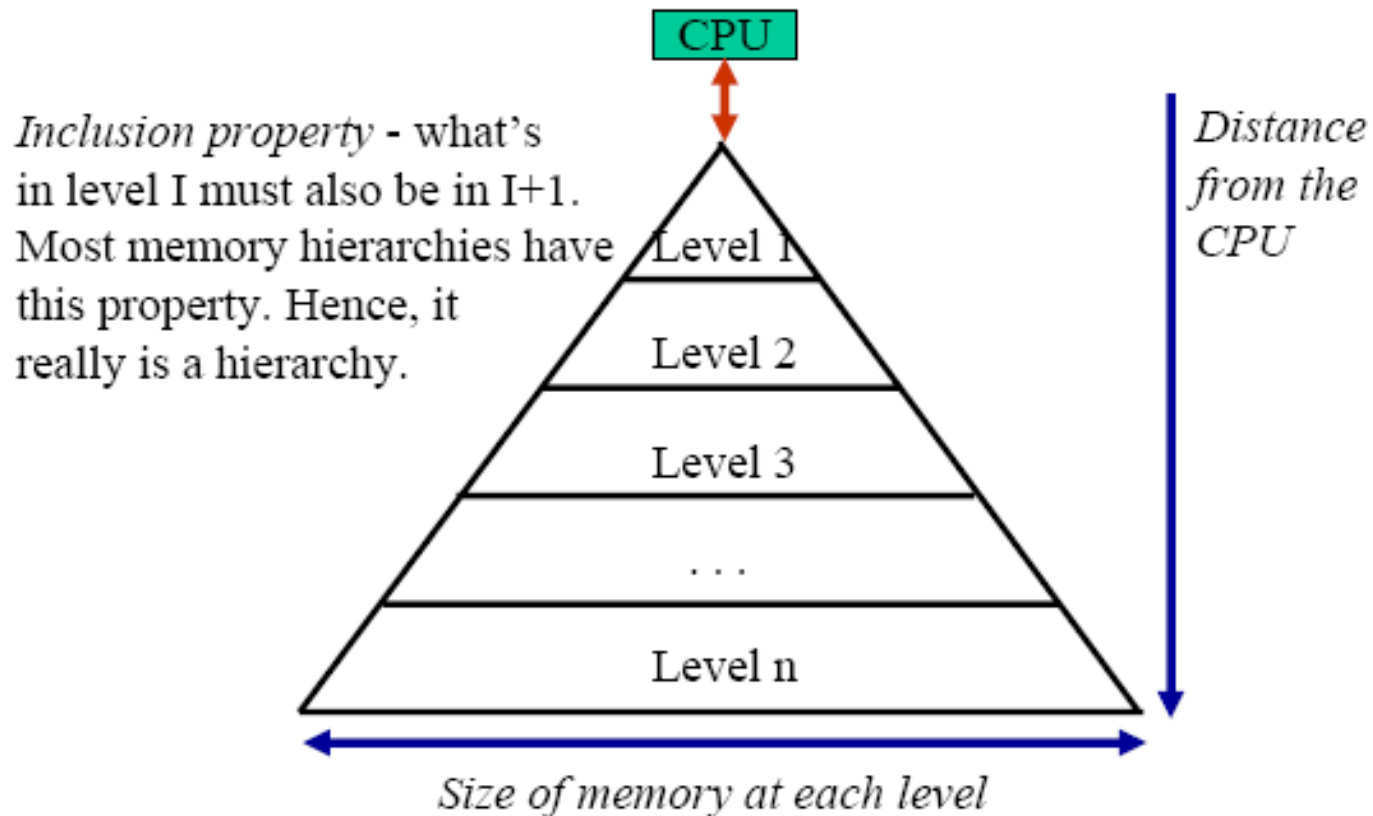
Longe da CPU – lento, grande e acesso menos freqüente aos dados



| | | | | | |
|-------------------------|-----|------|-------|-------|--------------------|
| Velocidade (ns): | 0,1 | 1 | 2-5 | 10-20 | 10.000.000 (10 ms) |
| Tamanho (bytes): | 100 | 16 K | 512 K | 256 M | G |

Hierarquia de Memória

- **Comportamento como uma Pirâmide**



Hierarquia de Memória

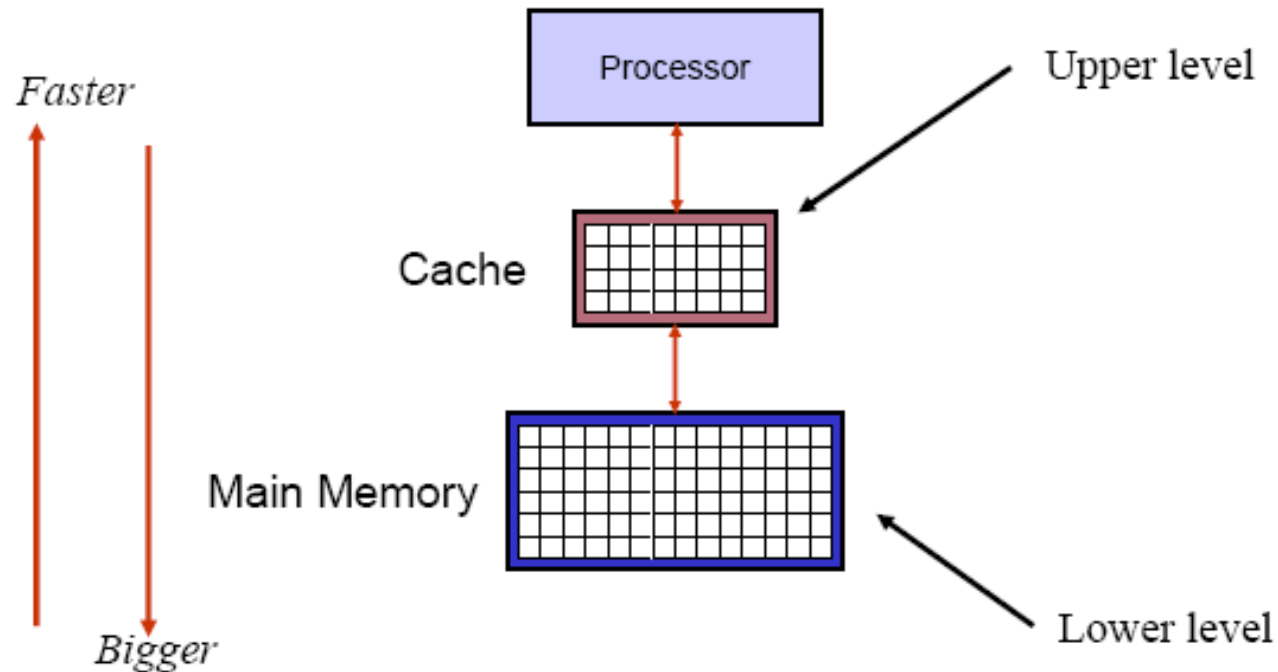
Como a hierarquia é gerenciada?

- **Registradores <-> memória**
 - pelo compilador
- **Cache <-> memória principal**
 - pelo hardware
- **Memória principal <-> disco**
 - pelo hardware e pelo sistema operacional (memória virtual)
 - pelo programador (arquivos)

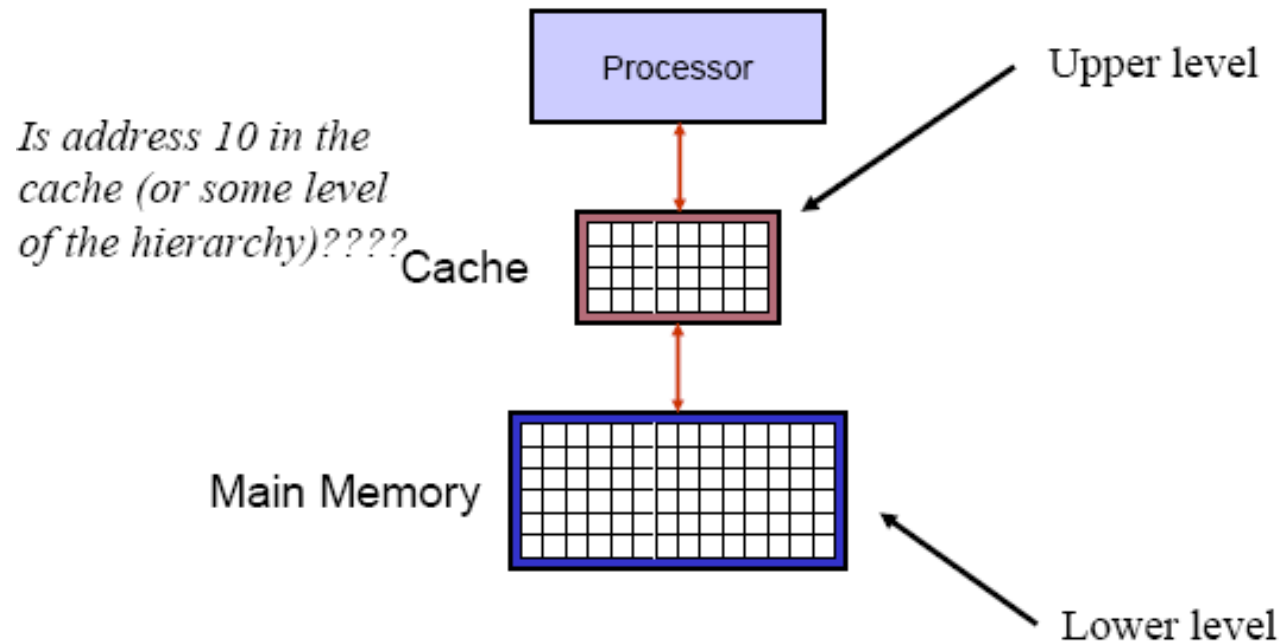
Hit e Miss - Acerto e Falta

- **Hit**: dado aparece em algum bloco no nível superior (junto ao processador)
 - **Hit Ratio**: razão de acertos dos 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**: falta do dado: precisa ser buscado de um bloco no nível inferior
 - **Miss Ratio** **razão da falta** = $1 - (\text{Hit Ratio})$
 - **Miss Penalty**: tempo gasto para substituir um bloco no nível superior + tempo para fornecer o bloco ao processador
- **Hit Time** \ll **Miss Penalty**

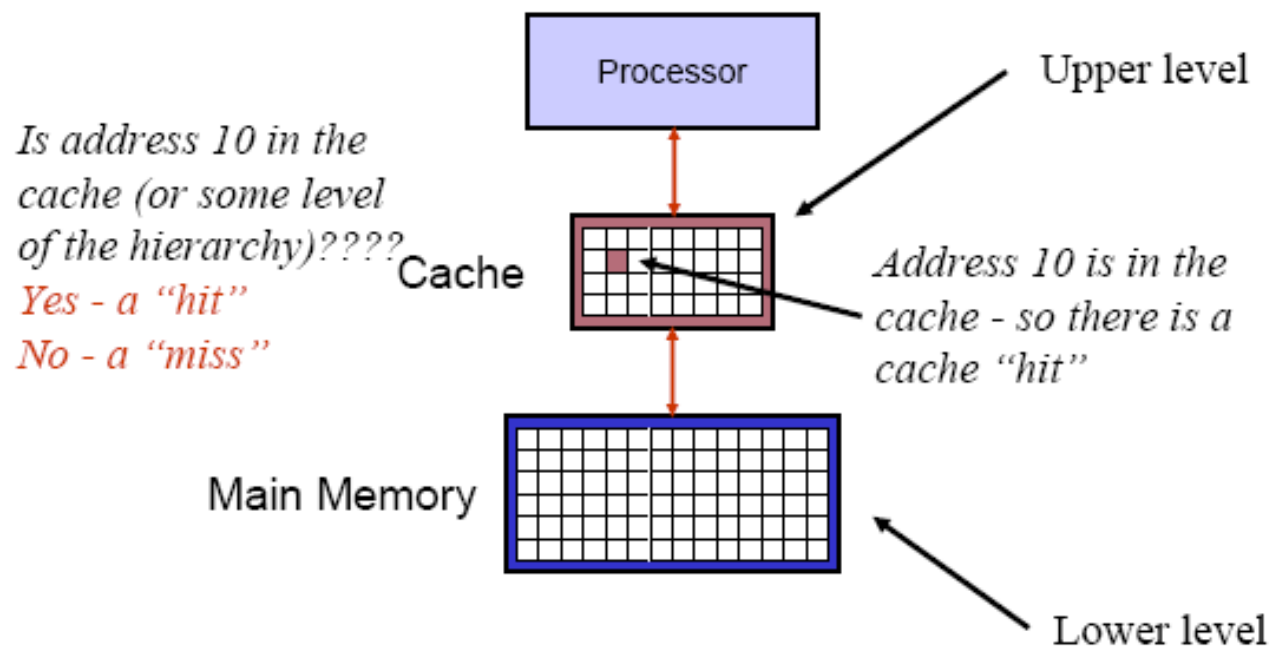
Hit e Miss



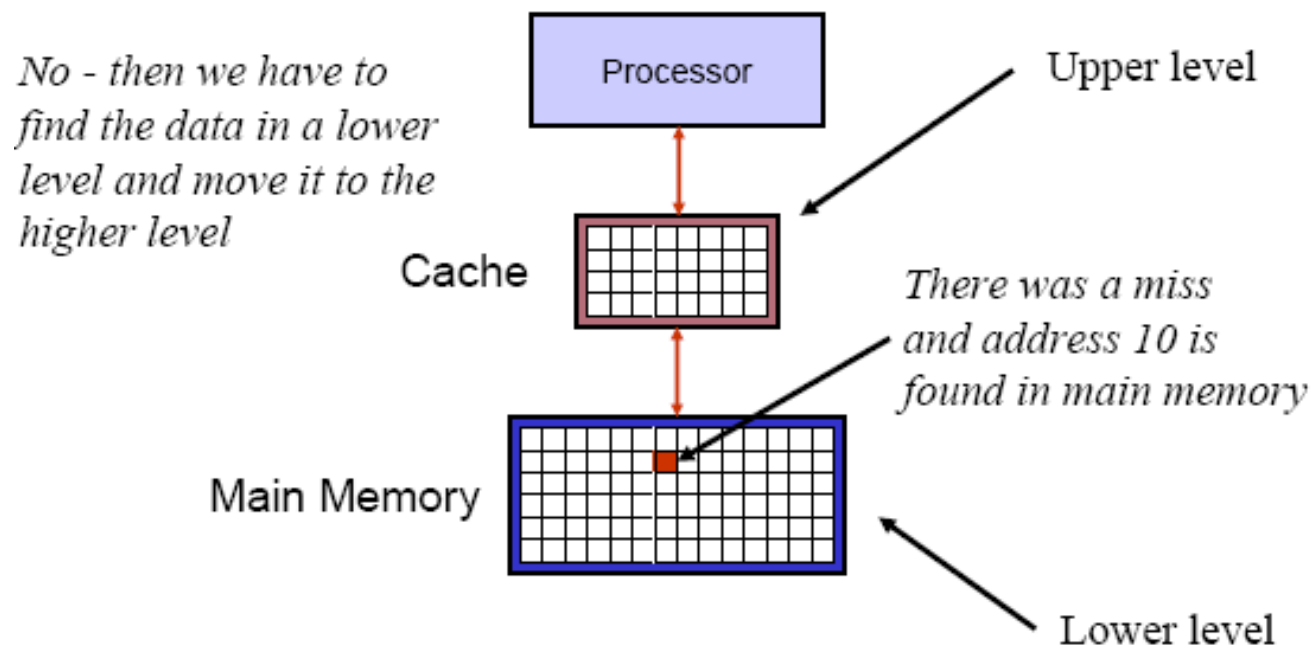
Hit e Miss



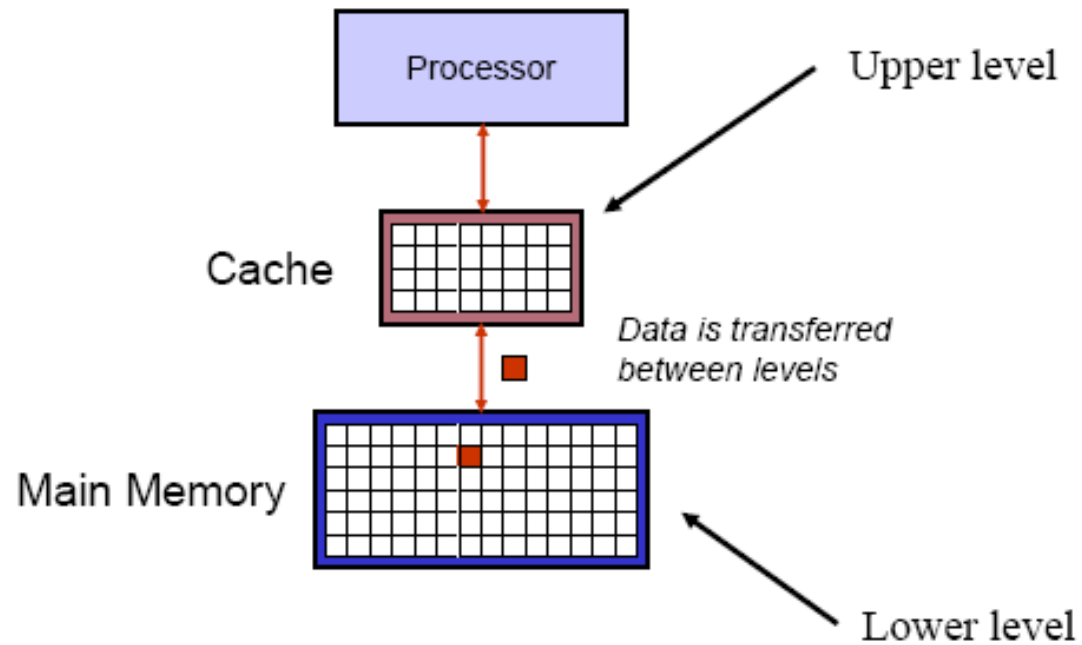
Hit e Miss



Hit e Miss

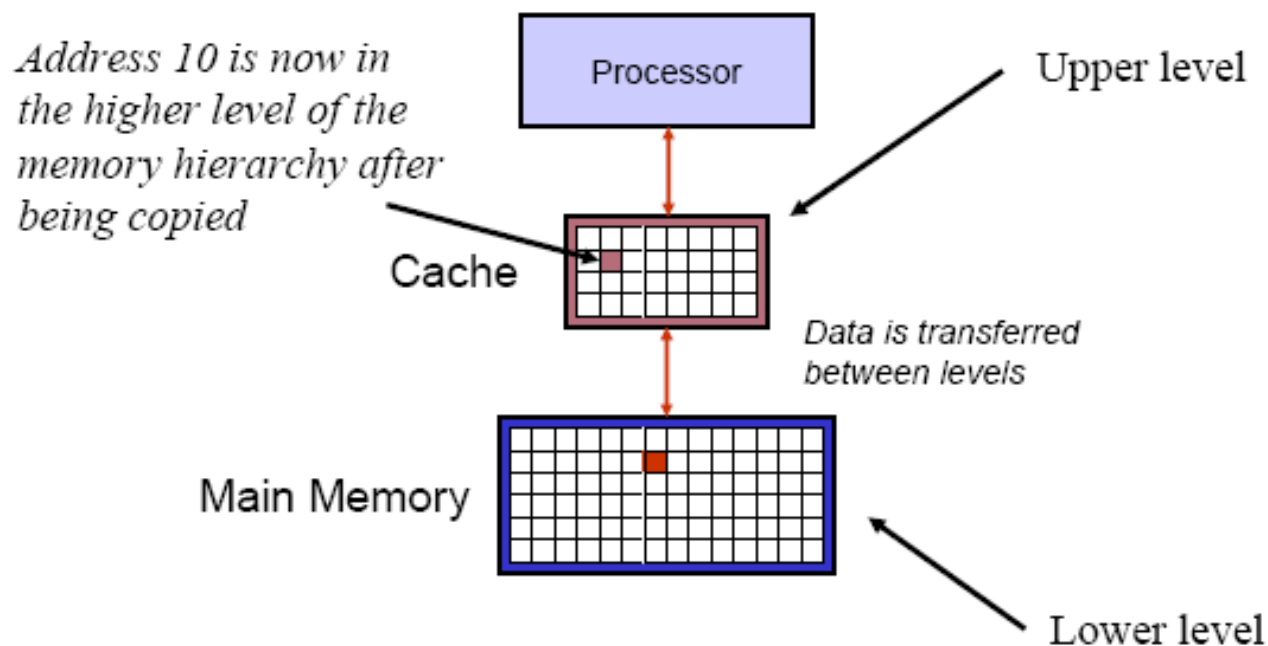


Hit e Miss



Hit e Miss

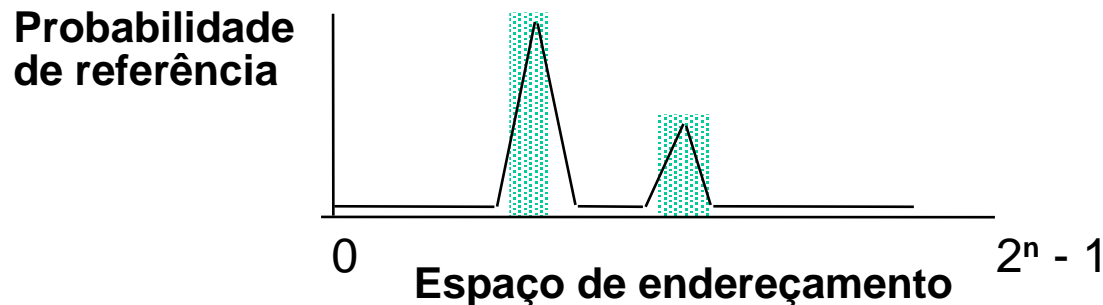
Bloco – menor unidade de transferência



3. Princípio de Localidade

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

Princípio de Localidade

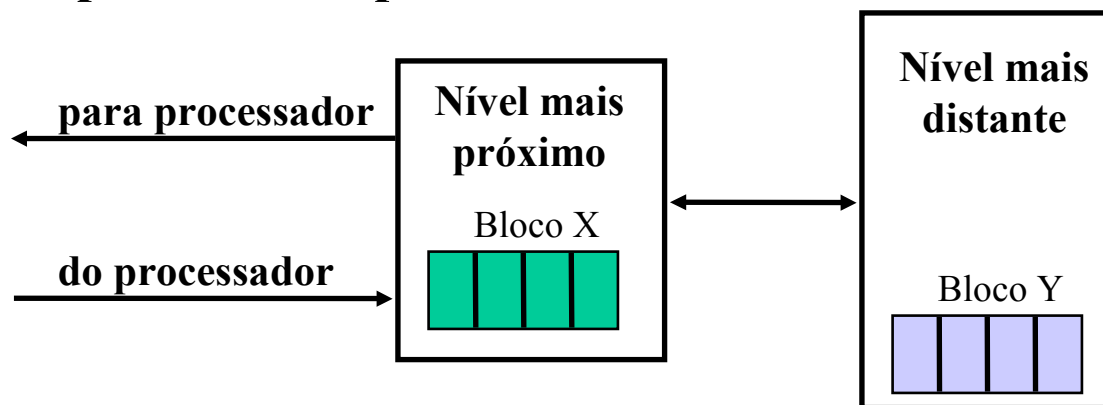
Como explorar o princípio de localidade numa hierarquia de memória?

- **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



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

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

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

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**

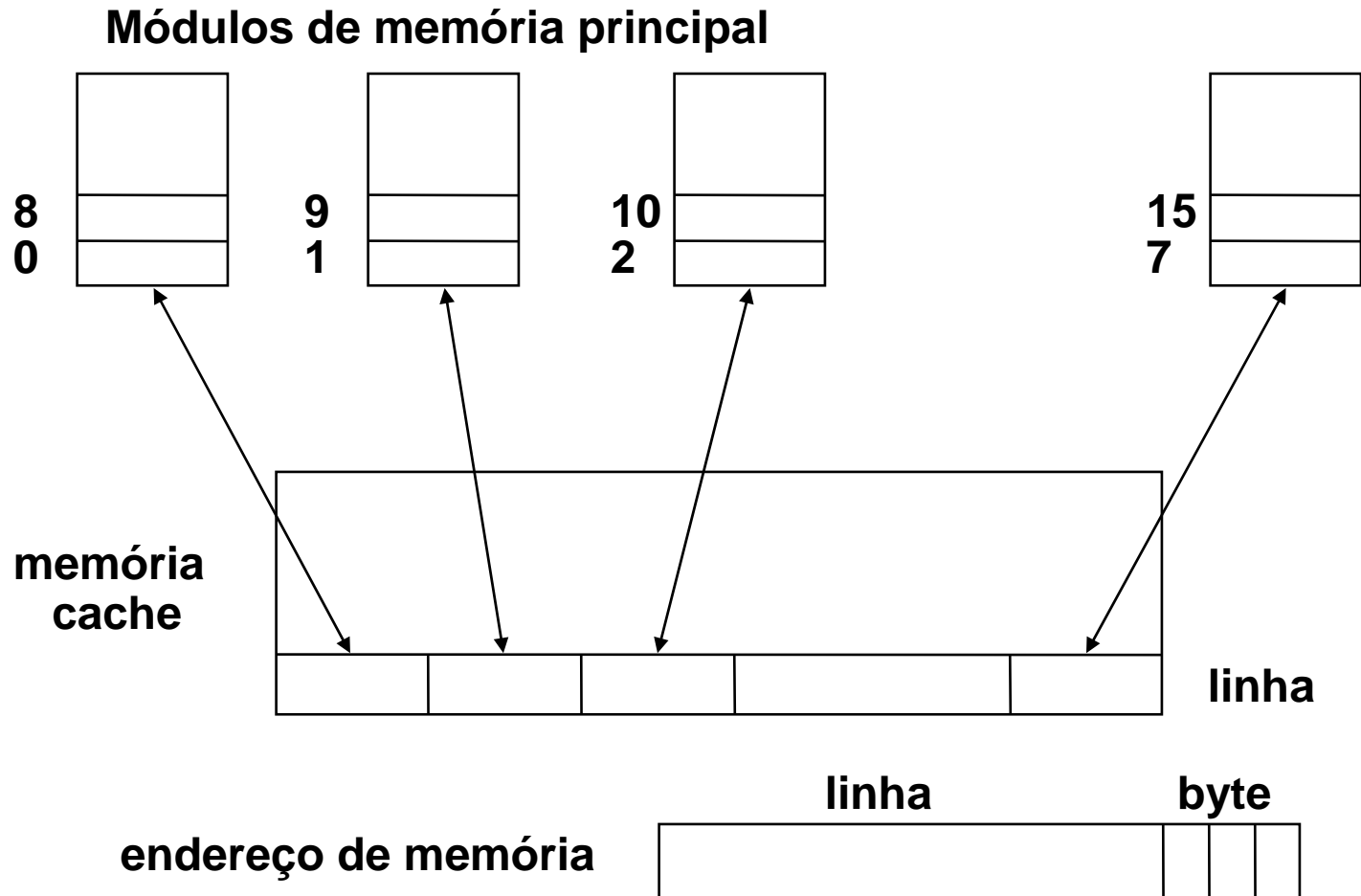
M = número de módulos da memória principal

T_c = tempo de acesso à cache

T_m = tempo de acesso à memória principal

Ideal: $M T_c = T_m$

Localidade Espacial



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

$$T_c = 1 \text{ ns}$$

$$T_m = 20 \text{ ns}$$

$$N = 10$$

$$N = 100$$



então $T_{ce} = 1,1 \text{ ns}$

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

4. Impacto no Desempenho

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** **0.95** **0.99** **1.0**



então $T_{ce} =$ **4 ns** **2 ns** **1.2 ns** **1 ns**

Impacto no Desempenho

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

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

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

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

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

| | |
|--------------|-----|
| 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

5. Organizações de Memória Cache

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
- 3 Estratégias de organização (mapeamento) da cache
 - mapeamento completamente associativo
 - mapeamento direto
 - mapeamento set-associativo

FIM