

# Hierarquia da Memória: Conceitos Fundamentais e Desempenho

Arquitetura de Computadores  
Lic. em Engenharia Informática

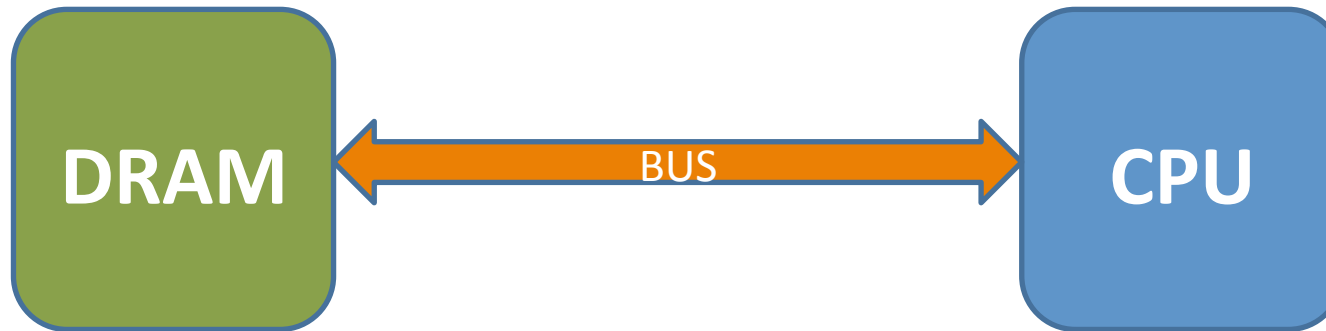
João Luís Sobral

# Hierarquia da Memória:

## Conceitos Fundamentais e Desempenho

|                                   |   |
|-----------------------------------|---|
| <b>Conteúdos</b>                  | 2.1 – Conceitos Fundamentais  |
|                                   | 2.2 – Hiato Processador-Memória   |
|                                   | 2.3 – Localidade  |
|                                   | 2.4 - Desempenho  |
| <b>Resultados de Aprendizagem</b> | R2.1 – Descrever e justificar a necessidade e oportunidade da hierarquia de memória |
|                                   | R2.2 – Quantificar o impacto da hierarquia da memória no desempenho da máquina      |

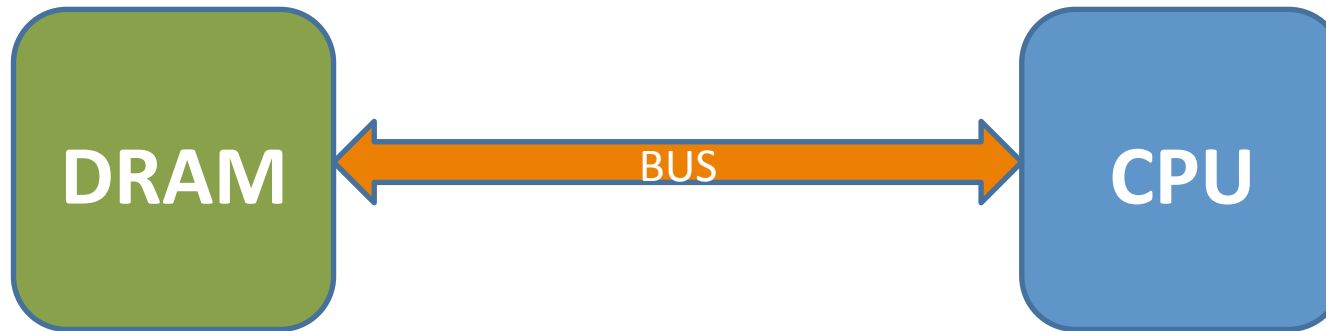
# Hiato Processador-Memória



Para cada instrução:

- 1.Ler instrução
- 2.Ler operando
- 3.Escrever Resultado

# Hiato Processador-Memória



Suponhamos um processador a executar um programa que consiste numa longa sequência de instruções inteiras:

```
addl reg, [Mem]
```

Se a instrução tiver 6 bytes de tamanho e cada inteiro 4 bytes a execução destas instruções implica um movimento de  $6 + 2 * 4 = 16$  bytes.

Se frequência = 2.5 GHz e o CPI=1 então são executadas  $2.5 * 10^9$  inst/seg  
A largura de banda necessária para manter o processador alimentado é de:

$$2.5 * 10^9 * 16 = \mathbf{40\ GB/s}$$

# Hiato Processador-Memória



| Standard name<br>(single channel) | Peak transfer<br>rate |
|-----------------------------------|-----------------------|
| DDR2-400                          | 3200 MB/s             |
| DDR2-800                          | 6400 MB/s             |
| DDR2-1066                         | 8533 MB/s             |
| DDR3-1066                         | 8533 MB/s             |
| DDR3-1600                         | 12800 MB/s            |

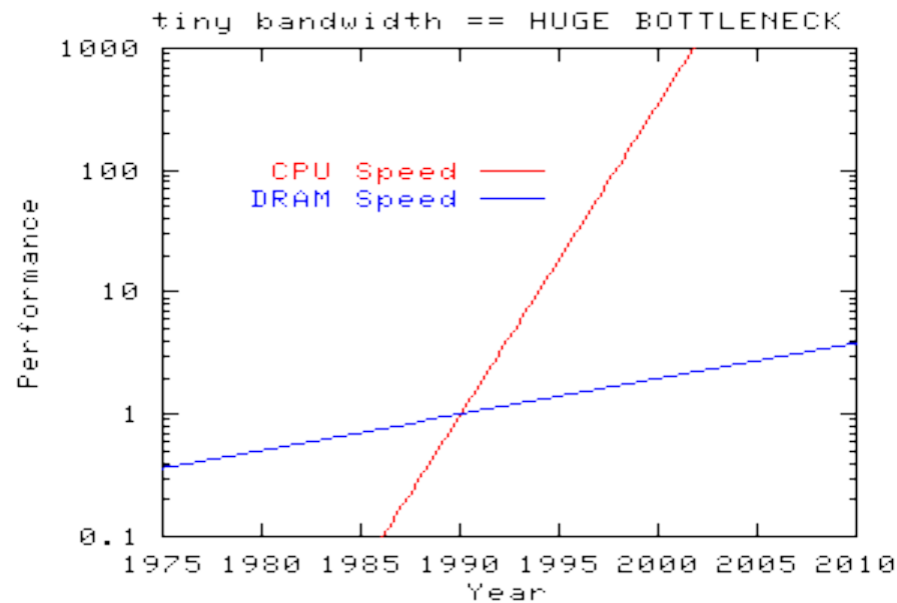
Largura de banda exigida neste exemplo:  $2.5 \cdot 10^9 \cdot 16 = \mathbf{40\ GB/s}$

Hiato processador-memória:

“A memória é incapaz de alimentar o processador com instruções e dados a uma taxa suficiente para o manter constantemente ocupado”

# Hiato Processador-Memória

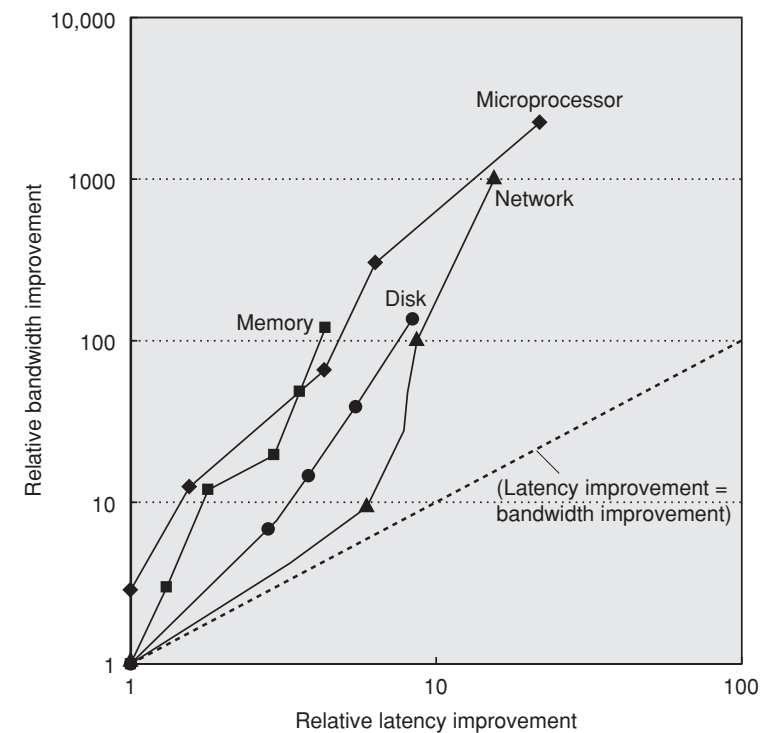
- As diferentes taxas de aumento do desempenho destes dois componentes levam a um aumento do hiato Processador-Memória (*“the memory gap”*) com o tempo
- O hiato processador-memória é um dos principais obstáculos à melhoria do desempenho dos sistemas de computação



The STREAM benchmark  
<http://www.cs.virginia.edu/stream/ref.html>

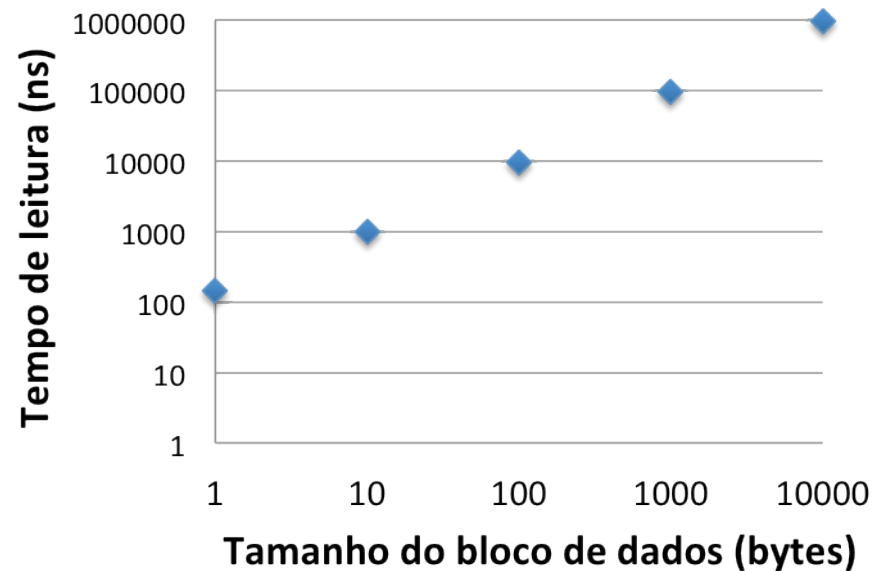
# Soluções: Hierarquia de Memória

- As memórias mais rápidas têm um custo por *byte* armazenado mais elevado
  - SRAM: 1000€ / GB; débito 100 GB/s
  - DRAM: 10€ / GB; débito 10 GB/s
  - Disco rígido: 0,10 € / GB; débito 100 MB/s
- Latência *versus* largura de banda (débito):
  - A evolução tem privilegiado a melhoria da largura de banda, em detrimento da latência:
    - Disco rígido WD 1TB WD10EARX;
      - Latência: 9 ms;
      - Débito 110MB/s
    - RAM DDR3-1333
      - Latência: 50 ns;
      - Débito: 10600 MB/s



# Soluções: Hierarquia de Memória

- Latência versus débito
  - RAM DDR3-1333
    - Latência: 50 ns;
    - Débito: 10600 MB/s = 0,95ns/byte (limite teórico)





# Hierarquia da Memória: Localidade

O **princípio da localidade**, exibido pela maior parte dos programas no acesso à memória, permite acelerar os acessos à mesma através de uma hierarquia

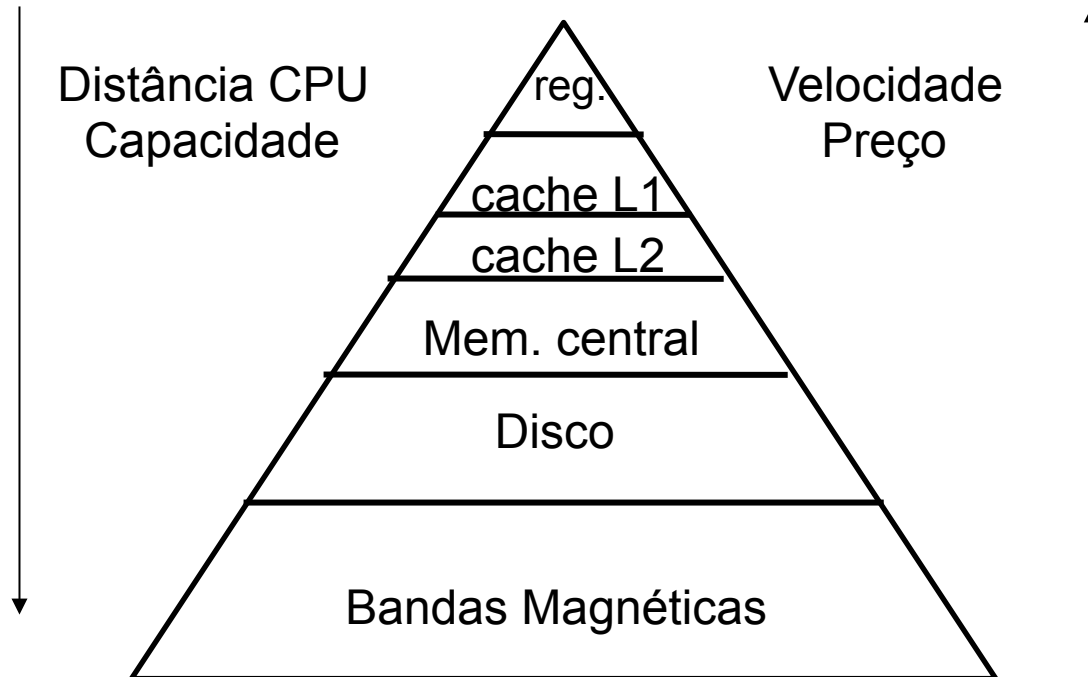
*“Os programas tendem a aceder a uma porção limitada de memória num dado período de tempo”*

O **princípio da localidade** permite utilizar memória mais rápida para armazenar a informação usada mais frequentemente/recentemente

Também permite tirar partido da largura de banda, uma vez que a informação transferida entre diferentes níveis da hierarquia é efectuada por blocos

# Soluções: Hierarquia de Memória

- Dotar a máquina de vários níveis de memória, tão mais rápidos (e mais caros por *byte* e com menor capacidade) quanto mais perto se encontram do processador.
- Cada nível contém uma cópia do código e dados mais usados em cada instante.



# Hierarquia da Memória: Localidade Temporal

**Localidade Temporal** – um elemento de memória acedido pelo CPU será, com grande probabilidade, acedido de novo num futuro próximo.

**Exemplos:** tanto as instruções dentro dos ciclos, como as variáveis usadas como contadores de ciclos, são acedidas repetidamente em curtos intervalos de tempo.

**Consequência** – a 1ª vez que um elemento de memória é acedido deve ser lido do nível mais baixo (por exemplo, da memória central).

Mas da 2ª vez que é acedido existem grandes hipóteses que se encontre na *cache*, evitando-se o tempo de leitura da memória central.

# Hierarquia da Memória: Localidade Espacial

**Localidade Espacial** – se um elemento de memória é acessado pelo CPU, então elementos com endereços na proximidade serão, com grande probabilidade, acessados num futuro próximo.

**Exemplos:** as instruções do programa são normalmente são acessadas em sequência, assim como, na maior parte dos programas, os elementos de vectores/matrizes.

**Consequência** – a 1ª vez que um elemento de memória é acessado, deve ser lido do nível mais baixo (por exemplo, memória central) não apenas esse elemento, mas sim um **bloco** de elementos com endereços na sua vizinhança.

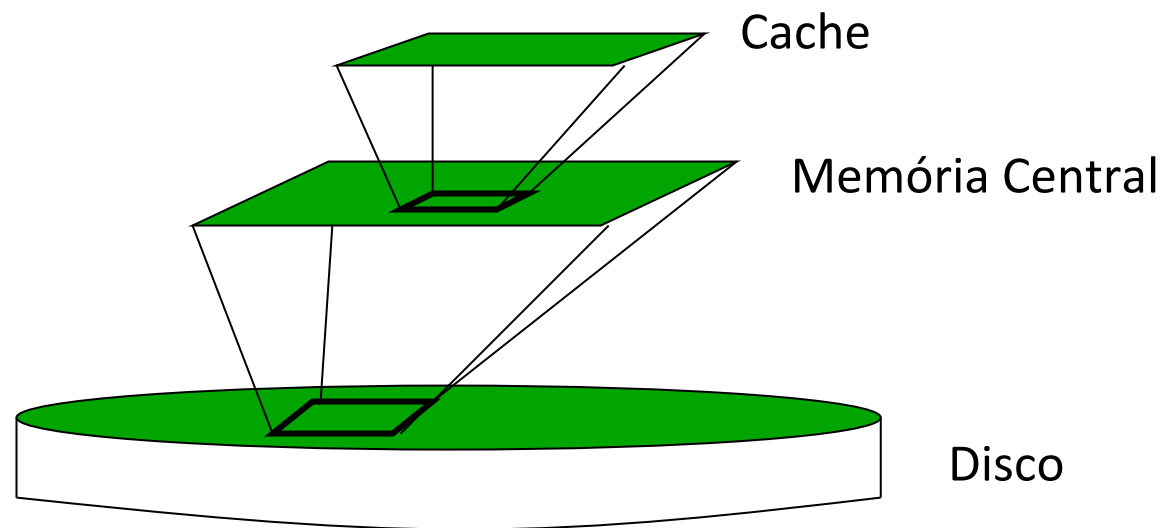
Se o processador, nos próximos ciclos, acede a um endereço vizinho do anterior (ex.: próxima instrução ou próximo elemento de um vector) aumenta a probabilidade de esta estar na *cache*.

# Hierarquia de Memória: Inclusão

Os dados contidos num nível mais próximo do processador são um sub-conjunto dos dados contidos no nível anterior.

O nível mais baixo contém a totalidade dos dados.

Os dados são copiados entre níveis em *blocos*



# Hierarquia de Memória: Terminologia

**Linha** – a cache está dividida em linhas. Cada linha tem o seu endereço (índice) e tem a capacidade de um bloco

**Bloco** – Quantidade de informação que é transferida de cada vez da memória central para a cache (ou entre níveis de cache). É igual à capacidade da linha.

**Hit** – Diz-se que ocorreu um *hit* quando o elemento de memória acedido pelo CPU se encontra na cache.

**Miss** – Diz-se que ocorreu um *miss* quando o elemento de memória acedido pelo CPU não se encontra na cache, sendo necessário lê-lo do nível inferior da hierarquia.

| Cache |     |
|-------|-----|
|       | 000 |
|       | 001 |
|       | 010 |
|       | 011 |
|       | 100 |
|       | 101 |
|       | 110 |
|       | 111 |

# Hierarquia de Memória: Terminologia

**Hit rate** – Percentagem de *hits* ocorridos relativamente ao total de acessos à memória.

**Miss rate** – Percentagem de *misses* ocorridos relativamente ao total de acessos à memória.  $Miss\ rate = (1 - hit\ rate)$

**Hit time** – Tempo necessário para aceder à *cache*, incluindo o tempo necessário para determinar se o elemento a que o CPU está a aceder se encontra ou não na cache.

**Miss penalty** – Tempo necessário para carregar um bloco da memória central (ou de um nível inferior) para a cache quando ocorre um *miss*.

# Hierarquia da memória - Desempenho

$$T_{exec} = \#I * CPI * T_{cc} \quad ou \quad T_{exec} = \#CC * T_{cc}$$

Como é que a hierarquia de memória influencia  $T_{exec}$ ?

Cada acesso à memória vai originar ciclos adicionais na execução do programa ( $\#CC_{MEM}$ ) devido aos *misses*:

$$T_{exec} = (\#CC_{CPU} + \#CC_{MEM}) * T_{cc}$$

Cada *miss* implica um aumento do  $\#CC$  em *misspenalty* ciclos, logo:

$$\#CC_{MEM} = \underbrace{n^{\circ} \text{ acessos à memória} * \text{miss rate}}_{n^{\circ} \text{ de miss}} * \text{miss penalty}$$



# Hierarquia da memória - Desempenho

$$T_{exec} = (\#CC_{CPU} + \#CC_{MEM}) * T_{cc}$$



$$\#CC = \#I * CPI$$

$$T_{exec} = \#I * (CPI_{CPU} + CPI_{MEM}) * T_{cc}$$

**CPI<sub>CPU</sub>** – nº de ciclos que o processador necessita, em média, para executar cada instrução;

O *hit time* considera-se incluído no CPI<sub>CPU</sub>

**CPI<sub>MEM</sub>** – nº de ciclos que o processador pára, em média, à espera de dados da memória, porque não encontrou estes dados na cache.

Estes são vulgarmente designados por ***memory stall cycles*** ou ***wait states***.

## Hierarquia da memória - Desempenho

$$CPI_{MEM} = \%acessosMem * missrate * misspenalty$$

Os acessos à memória devem-se a:

- acesso a **dados** (instruções de *Load* e *Store* do programa)
- busca de **instruções**

Como estes têm comportamentos diferentes usam-se diferentes percentagens:

- **Dados** – Apenas uma determinada percentagem de instruções acede à memória (%Mem). *missrate<sub>D</sub>* refere-se ao acesso a dados.
- **Instruções** – Todas as instruções são lidas da memória, logo a % de acesso à memória é de 100%. *missrate<sub>I</sub>* refere-se ao acesso às instruções.

*missrate<sub>I</sub>* é geralmente menor que *missrate<sub>D</sub>* devido à localidade espacial.

$$CPI_{MEM} = (missrate_I + \%Mem * missrate_D) * misspenalty$$

# Hierarquia da memória - Desempenho

Abreviando *missrate* por *mr* e *misspenalty* por *mp* temos

$$T_{exec} = \#I * (CPI_{CPU} + CPI_{MEM}) * T_{cc}$$

$$CPI_{MEM} = (mr_I + \%Mem * mr_D) * mp$$

substituindo

$$T_{exec} = \#I * [CPI_{CPU} + (mr_I + \%Mem * mr_D) * mp] * T_{cc}$$

**NOTA:** A *miss penalty* (*mp*) tem que ser expressa em ciclos do relógio.

# Hierarquia da memória - Desempenho

Considere uma máquina com uma *miss rate* de 4% para instruções, 5% para dados e uma *miss penalty* de 50 ciclos. Assuma ainda que 40% das instruções são *loads* ou *stores*, e que o  $CPI_{CPU}$  é 1. Qual o CPI total?

$$CPI = CPI_{CPU} + CPI_{MEM} = CPI_{CPU} + (mr_I + \%Mem * mr_D) * mp$$

$$CPI = 1 + (0.04 + 0.4 * 0.05) * 50 = 1 + 3 = 4$$

Se a frequência do relógio for de 1 GHz e o programa executar  $10^9$  instruções qual o tempo de execução?

$$T_{exec} = \#I * CPI * T_{cc} = 10^9 * 4 * \frac{1}{10^9} = 4s$$

# Hierarquia da memória - Desempenho

Considere um programa com as características apresentadas na tabela, a executar numa máquina com memória de tempo de acesso 0. Se a frequência do processador for 1 GHz, qual o CPI médio e o tempo de execução?

| Instrução     | Nº Instruções  | CPI |
|---------------|----------------|-----|
| Cálculo       | $3 \cdot 10^8$ | 1,1 |
| Acesso à Mem. | $6 \cdot 10^8$ | 2,5 |
| Salto         | $1 \cdot 10^8$ | 1,7 |
| <b>TOTAL:</b> | $10^9$         |     |

$$CPI = CPI_{CPU} + CPI_{MEM} = (3 \cdot 1.1 + 6 \cdot 2.5 + 1 \cdot 1.7) / 10 + 0 = 2$$

$$T_{exec} = \#I * CPI * T_{cc} = 10^9 * 2 * \frac{1}{10^9} = 2s$$

# Hierarquia da memória - Desempenho

Considere o mesmo programa e máquina do acetato anterior, mas agora com um tempo de acesso à memória de 10 ns (por palavra ou instrução). Suponha ainda que esta máquina não tem cache. Qual o CPI efectivo e  $T_{exec}$ ?

$$CPI = CPI_{CPU} + CPI_{MEM} = CPI_{CPU} + (mr_I + \%Mem * mr_D) * mp$$

Se a máquina não tem cache, então  $mr_I = mr_D = 100\%$ .

Da tabela tiramos que  $\%Mem = 60\%$ .

$mp$  expresso em ciclos do relógio é  $10/1 = 10$  ciclos ( $f=1$  GHz)

$$CPI = CPI_{CPU} + CPI_{MEM} = 2 + (1 + 0.6 * 1) * 10 = 2 + 16 = 18$$

$$T_{exec} = \#I * CPI * T_{cc} = 10^9 * 18 * \frac{1}{10^9} = 18s$$

# Hierarquia da memória - Desempenho

Considere agora que existe uma *cache* com linhas de 4 palavras; a *miss rate* de acesso às instruções é de 6% e de acesso aos dados é de 10%; o tempo de acesso à memória central é constituído por uma latência de 40 ns mais 10 ns por palavra. Qual o CPI médio e o tempo de execução?

$$mp = 40 + 10 * 4 = 80 \text{ ns} ; \text{ em ciclos } mp = 80 / 1 = 80 \text{ ciclos}$$

$$CPI = CPI_{CPU} + CPI_{MEM} = 2 + (0.06 + 0.6 * 0.1) * 80 = 2 + 9.6 = 11.6$$

$$T_{exec} = \#I * CPI * T_{cc} = 10^9 * 11.6 * \frac{1}{10^9} = 11.6s$$

# Hierarquia da memória - Desempenho

Suponha que a capacidade da *cache* é aumentada para o dobro, resultando numa *miss rate* de acesso às instruções de 3.2% e acesso aos dados de 8%. No entanto, o tempo de acesso à cache (*hit time*) também aumenta, resultando num  $CPI_{CPU}$  de 2.5 . Qual o CPI médio e o tempo de execução?

$$CPI = CPI_{CPU} + CPI_{MEM} = 2.5 + (0.032 + 0.6 * 0.08) * 80 = 2.5 + 6.4 = 8.9$$

$$T_{exec} = \#I * CPI * T_{cc} = 10^9 * 8.9 * \frac{1}{10^9} = 8.9s$$



# Hierarquia da memória - Desempenho

Para tirar maior partido da localidade espacial aumentou-se o número de palavras por linha de 4 para 8, reduzindo a *miss rate* de instruções para 1% e de dados para 6%. O tempo de acesso à memória central é composto por uma latência de 40 ns mais 10 ns por palavra. Qual o CPI médio e o tempo de execução?

$$mp = 40 + 10 * 8 = 120 \text{ ns ; em ciclos } mp = 120 / 1 = 120 \text{ ciclos}$$

$$CPI = CPI_{CPU} + CPI_{MEM} = 2.5 + (0.01 + 0.6 * 0.06) * 120 = 2.5 + 5.52 = 8.02$$

$$T_{exec} = \#I * CPI * T_{cc} = 10^9 * 8.02 * \frac{1}{10^9} = 8.02s$$

# Hierarquia da memória - Desempenho

Para reduzir a *miss penalty* a memória central foi substituída por outra com uma latência de 40 ns e 5 ns por palavra. Qual o CPI médio e o tempo de execução?

$$mp = 40 + 5 * 8 = 80 \text{ ns ; em ciclos } mp = 80 / 1 = 80 \text{ ciclos}$$

$$CPI = CPI_{CPU} + CPI_{MEM} = 2.5 + (0.01 + 0.6 * 0.06) * 80 = 2.5 + 3.68 = 6.18$$

$$T_{exec} = \#I * CPI * T_{cc} = 10^9 * 6.18 * \frac{1}{10^9} = 6.18s$$

# Hierarquia da memória - Desempenho

Finalmente o processador foi substituído por outro com uma frequência de 3 GHz, sem que a memória tenha sofrido qualquer alteração. Qual o CPI médio e o tempo de execução?

O ciclo do relógio é agora de 0.33 ns, logo  $mp = 80/0.33 = 240$  ciclos

$$CPI = CPI_{CPU} + CPI_{MEM} = 2.5 + (0.01 + 0.6 * 0.06) * 240 = 2.5 + 11.04 = 13.54$$

$$T_{exec} = \#I * CPI * T_{cc} = 10^9 * 13.54 * \frac{1}{3 * 10^9} = 4.513s$$