

Desempenho Cache

→ incluindo hazards e pipeline

↓ mr: associatividade, algoritmo
↓ penalidade: multilével

$$T_{\text{exec CPU}} = (\text{ciclos}_{\text{CPU}} + \text{ciclos}_{\text{paradas}}) \cdot T_{\text{CPU}}$$

↳ falta em cache

$$\text{ciclos}_{\text{paradas}} = \text{ciclos}_{\text{paradas}}^{(\text{leitura})} + \text{ciclos}_{\text{paradas}}^{(\text{escrita})}$$

→ miss rate

$$\rightarrow \text{ciclos}_{\text{paradas}}^{(\text{leitura})} = \frac{\text{leituras}}{\text{programa}} \cdot \text{mr}(\text{leitura}) \cdot \text{penalidade}(\text{leitura})$$

$$\rightarrow \text{Escrita: } \text{ciclos}_{\text{paradas}}^{(\text{escrita})} = \frac{\text{escritas}}{\text{programa}} \cdot \text{mr}(\text{escrita}) \cdot \text{penalidade}(\text{escrita}) \quad (*)$$

→ write buffer

- (*)
- write through: escreve na cache, escreve na memória
 - ↳ penalidade: acesso e escrita na memória principal
 - write back: só vai copiar para a memória quando for substituído
 - ↳ penalidade: tempo para copiar na memória principal
 - write allocate: deu write miss (***) e escreve nessa posição na cache
 - ↳ penalidade:

(***) quando você vai escrever na posição da cache e as tags não são iguais.

- ↳
- write allocate: escreve nessa posição na cache, destruindo o anterior.
 - ↳ Obs: vai evitar a escrita na memória (+ lento)
 - write-no-allocate: escreve nessa posição na memória.

$$\text{ex seja: } \text{ciclos}_{\text{paradas}} = \frac{\text{acessos memória}}{\text{programa}} \cdot \text{mr} \cdot \text{penalidade}$$

⇕

$$= \frac{\text{instruções}}{\text{programa}} \cdot \frac{\text{mr}}{\text{instrução}} \cdot \text{penalidade}$$

ex: Instruções

↳ mr: 2%

Dados

↳ mr: 4%

• CPI = 2

• penalidade = 100 ciclos

• % load store = 36%

$$1) I_{\text{ciclos miss}} = I \cdot 2\% \cdot 100 = 2I$$

$$D_{\text{ciclos miss}} = \underbrace{(I \cdot 36\%)}_{\text{peso das load/store}} \cdot 4\% \cdot 100 = 1,44I$$

$$2) \text{ciclos}_{\text{paradas}} = 3,44I$$

$$3) \div I \rightarrow \text{CPI} = 3,44$$

4) Tot CPI = $CPI_{base} + CPI_{paradas} = 2 + 3,44 = 5,44$

→ AMAT: average memory access time
 ↳ considera os hits e misses

Obs: usado para comparação de designs diferentes.

$$AMAT = T_{hit} + mrr \cdot \text{penalidade}$$

Ex:

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

$$\text{penalidade} = 20 \text{ ciclos}$$

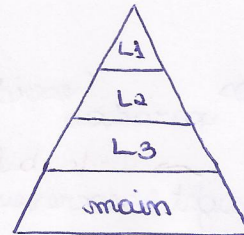
$$mrr = 0,05 = 5\%$$

$$\text{acesso} = 1 \text{ ciclo}$$

$$AMAT = 1 + 0,05 \cdot 20$$

$$AMAT = 2 \text{ ciclos}$$

Multilevel Caches



→ Motivo:

↳ retardar a busca na memória principal, adicionando caches maiores (penalidade menor que a memória) em níveis inferiores.

Ex: $CPI_{base} = 1$

$$\text{clock} = 46 \text{ Hz}$$

$$T = 0,25 \text{ ns}$$

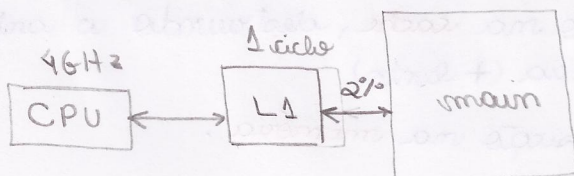
usar L1:

$$\text{penalidade (main)} = 100 \text{ ns}$$

$$I_{mrr} = 2\%$$

$$100 / 0,25$$

$$400 \text{ ciclos}$$



$$\text{Tot CPI} = CPI_{base} + CPI_{paradas}$$

$$= 1 + \frac{\text{ciclos paradas}}{I}$$

$$= 9$$

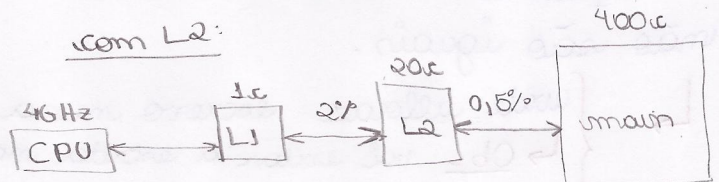
$$\rightarrow \text{ciclos paradas} = I \cdot mrr \cdot \text{penalidade}$$

$$= I \cdot 0,02 \cdot 400 = 8I$$

$$L2 \rightarrow \text{acesso} = 5 \text{ ns}$$

$$\text{miss rate global} = 0,5\%$$

com L2:



$$\text{Tot CPI} = CPI_{base} + CPI_{paradas}$$

$$= 1 + \frac{\text{ciclos paradas}}{I}$$

$$= 3,4$$

$$\rightarrow \text{ciclo paradas} = I \cdot (0,02 \cdot 20 + 0,005 \cdot 400)$$

$$= I \cdot 2,4$$



$$I [(2\% - 0,5\% \cdot 1 \cdot 20 + 0,5\% \cdot 400)]$$