

# Projeto de Caches

- Mapeamento de endereços (hashing)
- capacidade [bytes]
- tamanho de bloco [palavras]
- associatividade (mais hashing)
- três tipos de faltas
- tempo médio de acesso à memória

UFPR BCC CI212 2016-2— projeto de caches (i)

1

## Projeto de memórias cache (i)

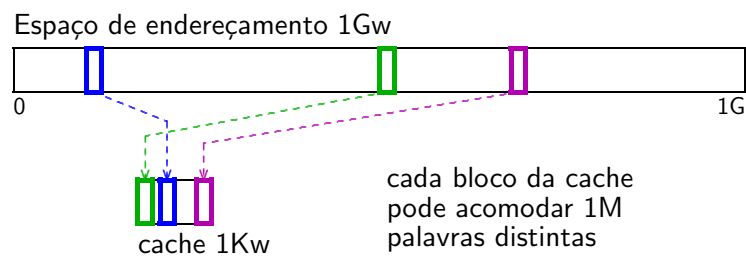
### Problema #1:

mapear espaço de endereçamento de 4Gbytes em 4Kbytes

necessita função que mapeia 1Gwords em 1Kwords

hashing

como implementa função em hardware, *eficientemente*?

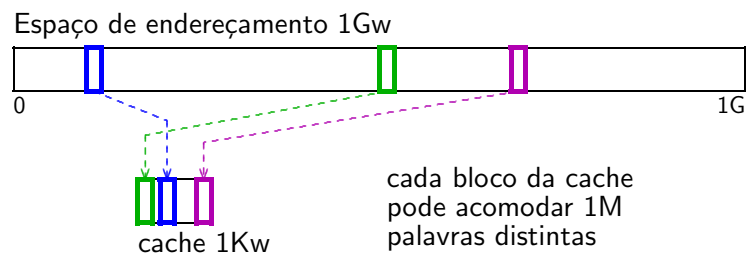


UFPR BCC CI212 2016-2— projeto de caches (i)

2

## Projeto de memórias cache (ii)

Problema #1: mapear espaço de 4Gbytes em 4Kbytes



UFPR BCC CI212 2016-2— projeto de caches (i)

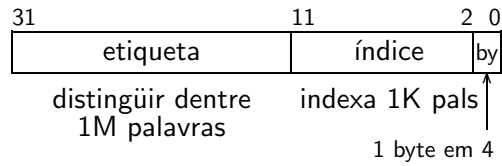
3

## Projeto de memórias cache (iii)

Mapeamento direto:

índice na cache é determinado por 10 bits do endereço

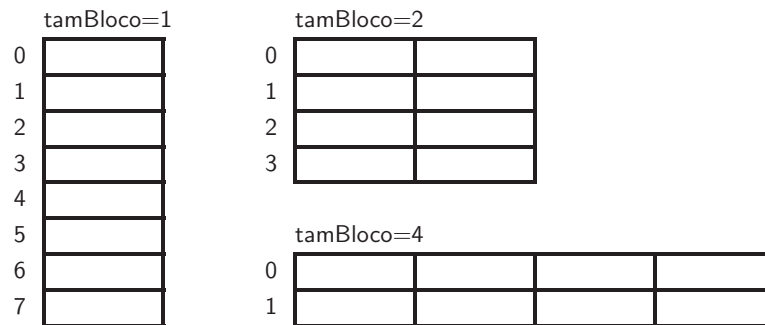
$$\text{pos} = \text{end\_WD} \% 1024$$



## Projeto de memórias cache (iv)

**Problema #2:** tirar vantagem da localidade espacial

necessita acomodar mais de uma palavra em cada bloco da cache

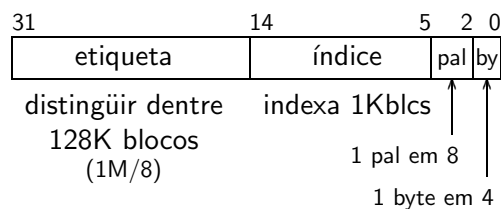


preencha as 3 caches com refs a palavras:

0 1 3 4 5 8 9 12 11 15 19 4 2 16 13 2 2 19 18 3

## Projeto de memórias cache (v)

**Problema #2:** tirar vantagem da localidade espacial

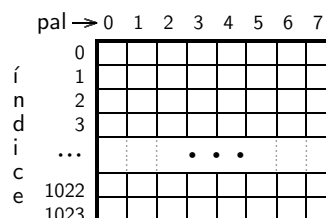


Mapeamento direto:

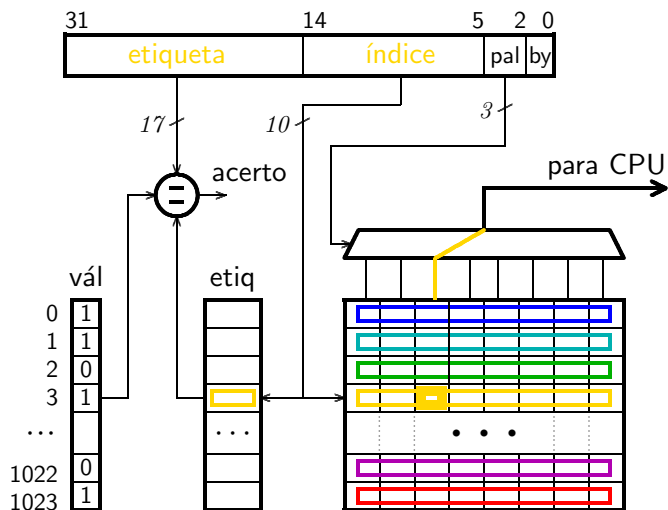
índice na cache é determinado por 10 bits do endereço

$$\text{pos} = \lfloor \text{end\_WD} / \text{tamBlc} \rfloor \% \text{númBlc}$$

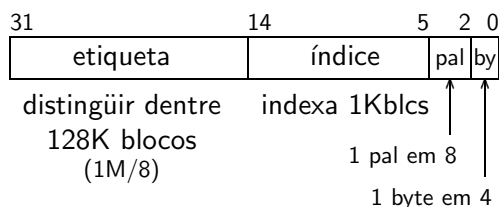
bloco acomoda 8 palavras de mesmo eti-q-índice



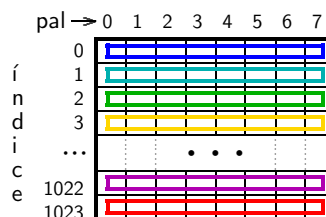
## Cache com Mapeamento Direto



## Mapeamento Direto e conflitos

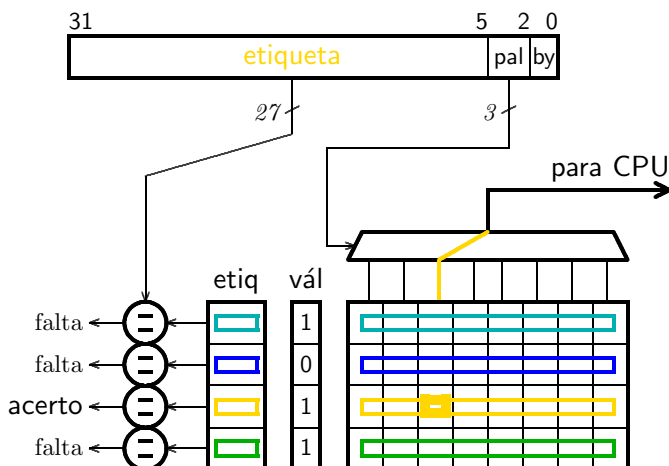


Mapeamento direto:  
 cada bloco da cache  
 pode armazenar palavras  
 de uma “cor” mas  
 memória contém 128K blocos  
 de cada cor (nesta cache 1K,8pal/bl)



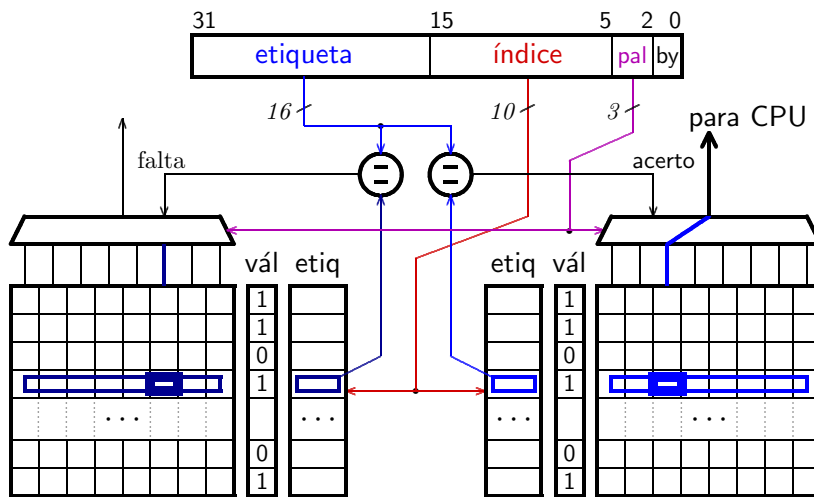
O que acontece se dois blocos da mesma cor são acessados em seqüência? índice igual, etiquetas distintas

## Mapeamento Associativo



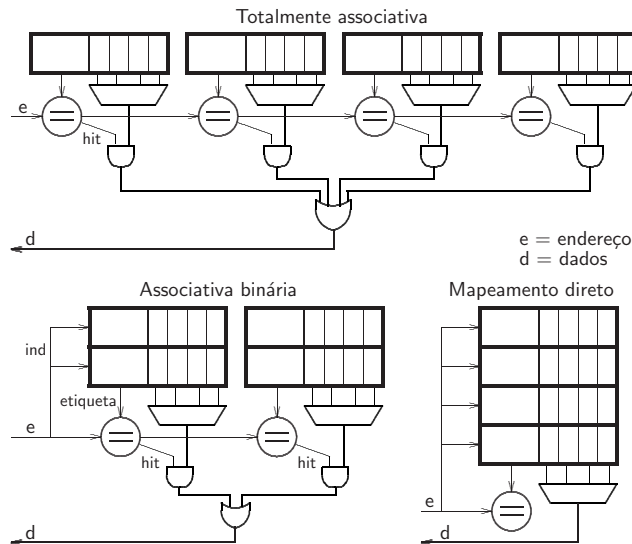
Mapeamento ideal e caro: um comparador em cada bloco da cache

# Mapeamento Híbrido: Associatividade por Conjuntos

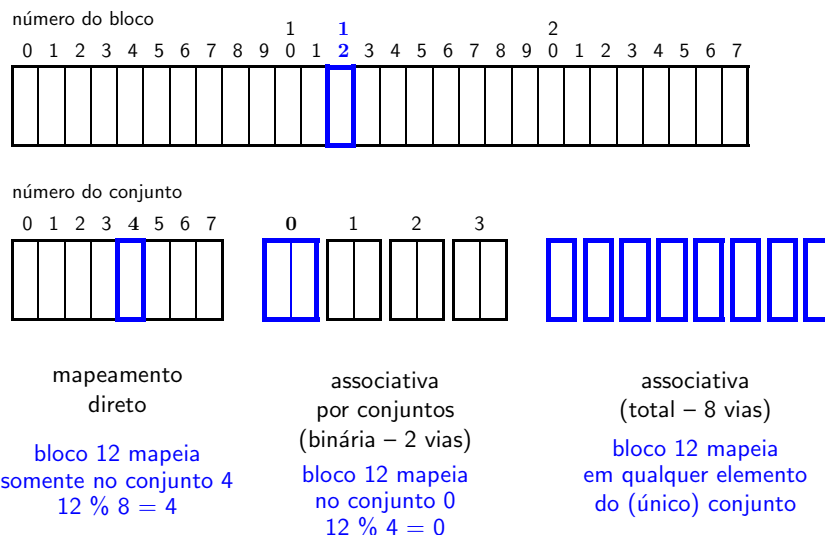


bloco de índice  $n$  pode estar em  $\forall$  elmtos do conjunto  $N$   $|N|=2$   
 cache com  $2^{|\text{índice}|}$  conjuntos de 2 elementos/conjunto

## Três Mapeamentos



## Alocação nos Blocos (i)



## Alocação nos Blocos (ii)

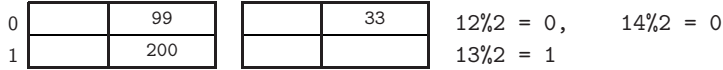
$M[12] := 99$ ;  $M[13] := 200$ ;  $M[14] := 33$

Associativa - 1 conjunto (4 elem/conj = 4 vias)

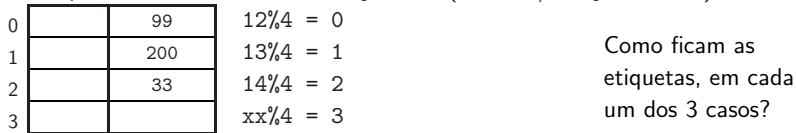


alocação:  $12\%1 = ?$

Associativa binária - 2 conjuntos (2 elem/conj = 2 vias)

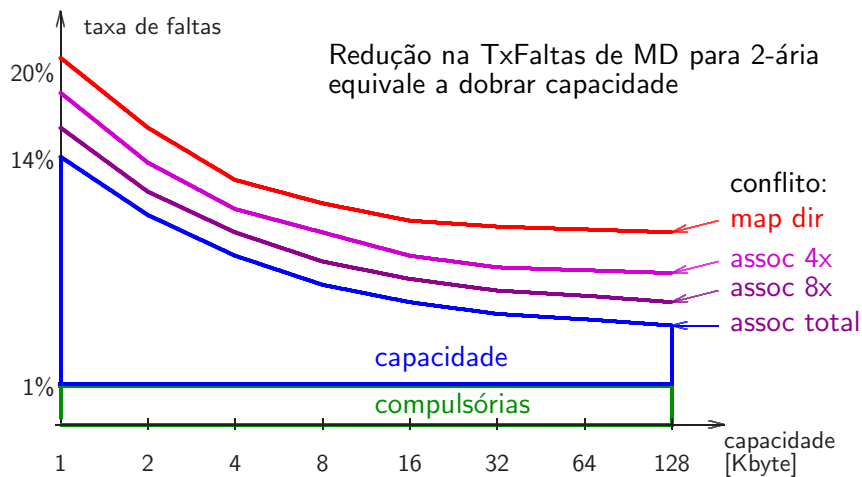


Mapeamento direto - 4 conjuntos (1 elem/conj = 1 via)



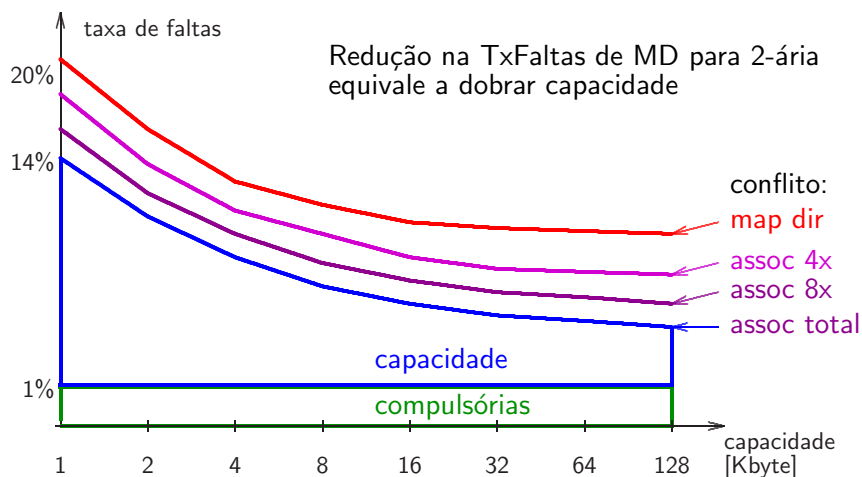
## Tipos de Faltas – falta por conflito

**problema:**  $N$  blocos da mesmo índice (côr) em uso simultâneo, mas cache só acomoda  $M < N$  blocos



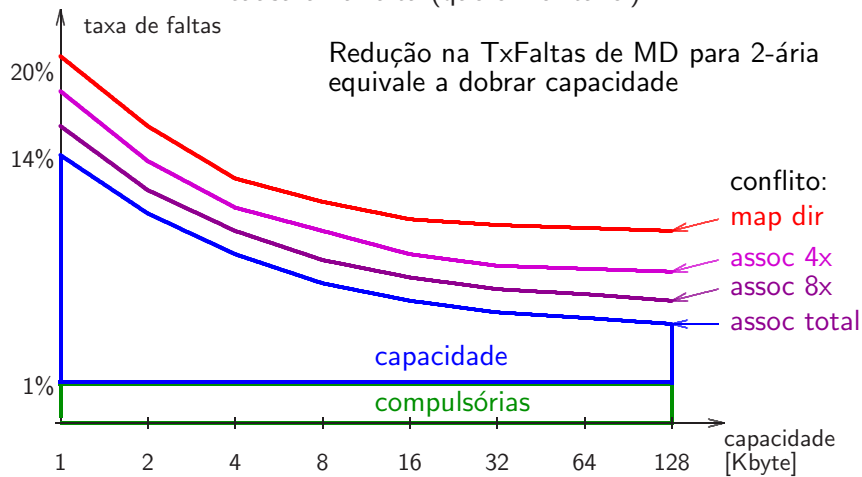
## Tipos de Faltas – falta por capacidade

**problema:** cache não tem capacidade para acomodar todos blocos acessados por programa



## Tipos de Faltas – falta compulsória

**problema:** primeiro acesso a um bloco **sempre** causa uma falta (que é inevitável)



UFPR BCC CI212 2016-2— projeto de caches (i)

16

## Tipos de Faltas – modelo qualitativo

### faltas por conflito

**problema:**  $N$  blocos da mesmo índice (côr) em uso simultâneo, mas cache só acomoda  $M < N$  blocos

**solução:** aumentar associatividade

### faltas por capacidade

**problema:** cache não tem capacidade para acomodar todos blocos acessados por programa

**solução:** aumentar tamanho da cache

### faltas compulsórias

**problema:** primeiro acesso a um bloco **sempre**

**solução:** busca antecipada por hardware ou por software

UFPR BCC CI212 2016-2— projeto de caches (i)

17

## Tipos de Faltas – 3 Cs de outro ângulo

- Que tipos de faltas ocorrem numa cache totalmente associativa de tamanho infinito?

faltas compulsórias porque deve trazer cada bloco para a cache

- Além destes, que tipos de faltas ocorrem numa cache totalmente associativa, de tamanho finito?

faltas por capacidade porque programa pode usar mais blocos do que a cache comporta

- Além destes, que tipos de faltas ocorrem numa cache associativa por conjunto ou com mapeamento direto?

faltas por conflito porque dois endereços  $\neq$ s mapeiam no mesmo bloco da cache (mesmo índice)

**Conflict misses** are misses that would not occur if the cache was fully associative and had LRU replacement [Jouppi90]

UFPR BCC CI212 2016-2— projeto de caches (i)

18

## Política de substituição de blocos

Em **caches associativas**, depois de uma falta de leitura, se não há blocos vazios, qual bloco deve ser expurgado da cache?

**bloco usado menos recentemente?**

Parece ser o melhor, mas pode ser difícil de implementar

**escolher aleatoriamente?**

É fácil de implementar, mas como é o desempenho?

TxFaltas Cache Assoc Binária

capac	LRU	aleatória
16KB	5.2%	5.7%
64KB	1.9%	2.0%
256KB	1.15%	1.17%

least recently used = LRU

Implementação:

LRU 2-way: bit aponta para LRU

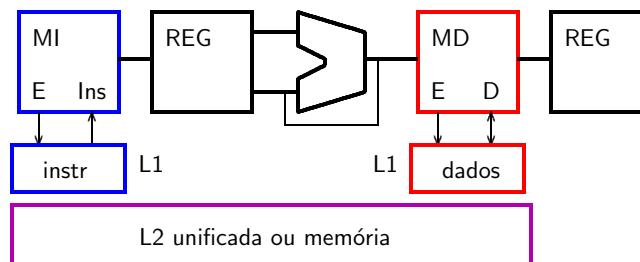
LRU 4-way: 3 bits, hierárquico

Aleatória: (a) somatório de alguns bits do endereço; (b) registrador de deslocamento com realimentação

## Caches e Pipelines (i)

**Tempo de Acerto** é crucial porque limita o ciclo da CPU

Se muito longo → mais de um ciclo para acessar a cache nos acertos



Há duas portas para memória para eliminar risco estrutural;

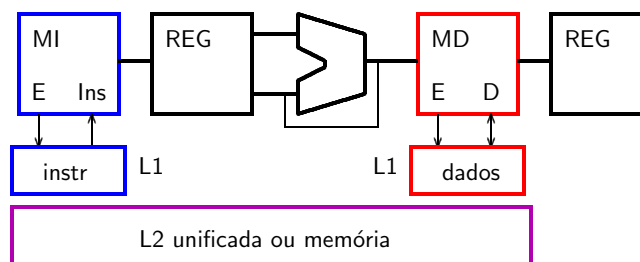
2 caches separadas, uma para instruções, uma para dados porque

localidade no acesso aos dados é diferente do acesso a instruções

## Caches e Pipelines (ii)

**Acertos:** se encontrou inst/dado na cache,

entrega para unidade funcional (RegInstr ou MEM) e prossegue



**Faltas:** segura execução até que memória entregue bloco faltante

**instrução:** segura busca, outras instruções prosseguem

**dado:** deve parar a busca até que memória entregue dado

## Desempenho do sistema com cache

$$\text{Equação do desempenho: } \frac{\text{segs}}{\text{prog}} = \frac{\text{instr}}{\text{prog}} \times \frac{\text{ciclos}}{\text{instr}} \times \frac{\text{segs}}{\text{ciclo}}$$

Até agora, CPI com tempo de acesso à memória constante:

$$\text{CPI}_{\text{real}} = \text{CPI}_{\text{ideal}} + \text{ciclos esperando pela memória}$$

$\text{CPI}_{\text{real}}$  depende do **Tempo Médio de Acesso à Memória**

$$\text{TMAM} = \text{TempoDeAcerto} + (\text{TaxaDeFaltas} \times \text{PenalidadePorFalta})$$

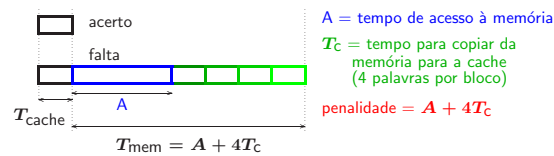
objetivo de projeto: **minimizar TMAM para dados e instruções**

**cuidado:** melhorar um termo pode piorar outros e aumentar TMAM

$$AMAT = \text{Average Memory Access Time}$$

## Desempenho do sistema com cache

$$\text{TMAM} = \text{TempoDeAcerto} + (\text{TaxaDeFaltas} \times \text{PenalidadePorFalta})$$



Compilador GCC com memória perfeita tem  $\text{CPI}=1.2$  (dependências), 11% das referências causam faltas que custam 10 ciclos cada.

$$\begin{aligned} \text{CPI}_{\text{real}} &= (\text{CPI}_{\text{ideal}} + \text{faltas} \times \text{penalidade}) \\ &= (1.2 + 0.11 \cdot 10) = 2.3 \end{aligned}$$

## Minimizar Tempo de Acerto

$$\text{TMAM} = \text{TempoDeAcerto} + (\text{TaxaDeFaltas} \times \text{PenalidadePorFalta})$$

- Tempo de acesso é proporcional à capacidade da memória  
→ caches pequenas
- Mapeamento Associativo é **mais lento** do que mapeamento direto  
acesso às etiquetas ; **comparação** ; **propagar através do MUX**
- Segmentar o acesso à cache (aumenta vazão, mantém latência  $\approx$ )  
endereçamento ; acesso aos dados ; transferência 2-3 estágios  
adiciona estágios ao pipeline da CPU:  
2-3 est em BUSCA + 2-3 em MEM



## Minimizar Taxa de Faltas

$$TMAM = \text{TempoDeAcerto} + (\text{TaxaDeFaltas} \times \text{PenalidadePorFalta})$$

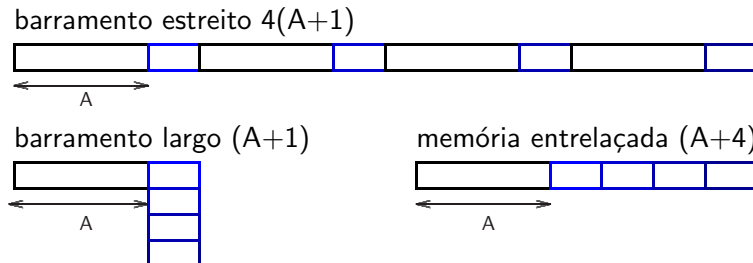
- Minimizar a taxa de faltas, ajustando para os 3Cs
  - \*  **aumentar capacidade**  tempo de acesso ↗  
 capacidade = num\_blocos × tam\_bloco × tam\_palavra
  - \*  **aumentar associatividade**  tempo de acesso ↗  
 capacidade = (n\_blocos × associativ) × t\_bloco × t\_palavra
  - \*  **busca antecipada**  por hardware ou por software  
 cuidado para não desalojar blocos úteis poluição  
 traz blocos e armazena em buffer específico

$$\begin{aligned} \text{capacidade MD} &= \text{num\_blocos} \times \text{tam\_bloco} \times \text{tam\_palavra} \\ [\text{byte}] &= [\text{bloco}] \times [\text{pal/bloco}] \times [\text{byte/pal}] \end{aligned}$$

## Minimizar Penalidade por Falta (i)

$$TMAM = \text{TempoDeAcerto} + (\text{TaxaDeFaltas} \times \text{PenalidadePorFalta})$$

- Minimizar o tempo de transferência com barramento mais eficiente
  - barramento largo (e caro);  
cache organizada em bancos;  
barramento CPU-cache com transações;
  - cache de vítimas: cache associativa  
pequena (4-8 blcs) mantém blocos  
expurgados da L1  
blocos na CdV tem 2a chance



## Minimizar Penalidade por Falta (ii)

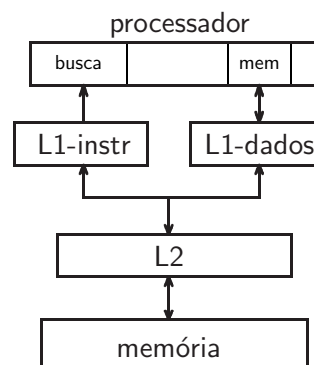
$$TMAM = \text{TempoDeAcerto} + (\text{TaxaDeFaltas} \times \text{PenalidadePorFalta})$$

Minimizar tempo de acerto  
no nível inferior (L2 ou memória)

$$Tx_{\text{FaltasLocal}} L_n = \frac{\text{faltas-}L_n}{\text{refs-}L_n}$$

$$Tx_{\text{FaltasGlobal}} = \frac{\text{faltas-}L_{1,2}}{\text{refs-CPU}}$$

Taxa de faltas na L2 é alta porque  
L1 filtra referências com boa localidade,  
sobram para L2 refs com localidade ruim



$$T_{\text{mam}} = T_{L1} + F_{L1} \cdot [T_{L2} + F_{L2} \cdot T_m]$$

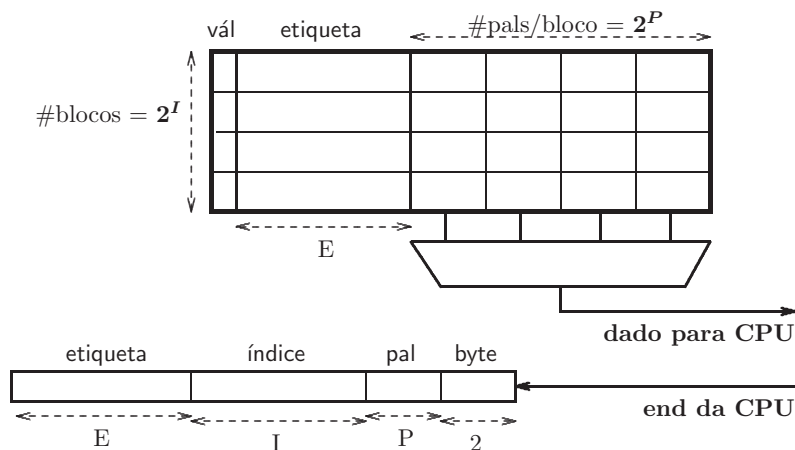
## Minimizar Penalidade por Falta (iii)

$$TMAM = \text{TempoDeAcerto} + (\text{TaxaDeFaltas} \times \text{PenalidadePorFalta})$$

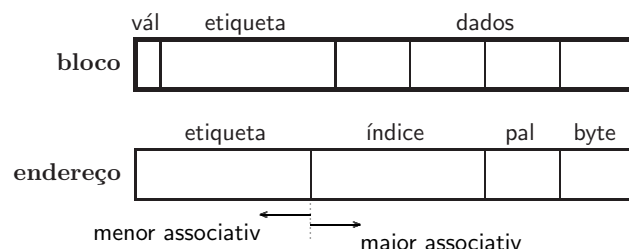
- **Minimizar tempo de carga:** preenchimento do bloco é demorado
- ...mas não precisa esperar até que **todo** bloco seja preenchido
- **Early restart** – assim que palavra requisitada chegar da memória, entrega ao processador, que continua a executar (instruções)
- **Critical word first** – busca palavra requisitada primeiro e a entrega ao processador assim que chegar da memória (CPU continua)
  - \* processador requisita **p2**: 

p0	p1	<b>p2</b>	p3
----	----	-----------	----
  - \* **p2** é entregue, e bloco é preenchido com  $p2 \rightarrow p3 \rightarrow p0 \rightarrow p1$
  - \* localidade espacial indica que próximo acesso será na palavra seguinte, que pode não ter chegado ainda

### Organização – Detalhes



### Organização – Detalhes



$\text{capacidade} = \text{num\_blocos} \times \text{tam\_bloco} \times \text{tam\_palavra} \times \text{associativ}$   
 $|\text{índice}| = \log_2 \text{num\_blocos}$       # conjuntos, altura da matriz  
 $|\text{pal}| = \log_2 \text{tam\_bloco}$       # palavras/dados, largura  
 $|\text{byte}| = \log_2 \text{tam\_palavra}$   
 $|\text{etiqueta}| = \text{tam\_ender} - ( (|\text{índice}|/\text{associativ}) + |\text{pal}| + |\text{byte}| )$   
 $\text{associatividade} = |\text{conjunto}|$

## Resumo

- Hashing para mapear espaço grande em armazenador pequeno se largura aumenta ( $\text{tamBloco} > 1$ ) então altura diminui
- Três formas de hashing:
  - \* mapeamento direto  $\rightarrow \text{posição} = \text{ender} \% \text{num\_blocos}$
  - \* totalmente associativo  $\rightarrow \text{posição} = \forall$  no conjunto (índice)
  - \* associativo por conjuntos  $\rightarrow \text{posição} = \text{ender} \% \text{num\_conjs}$
- Três categorias de faltas:
  - \* por capacidade  $\rightarrow$  cache maior
  - \* por conflitos  $\rightarrow$  associatividade
  - \* compulsórias  $\rightarrow$  busca antecipada (?)
- $\text{TMAM} = \text{TempoDeAcerto} + (\text{TaxaDeFaltas} \times \text{PenalidadePorFalta})$   
melhorar um termo geralmente piora outro/s

## Exercícios

- 1) Projete uma cache com 512Kbytes de capacidade, com blocos de 32 bytes, e associatividades:  
(a) mapeamento direto, (b) associativ binária, (c) associativ 8-ária.  
Indique claramente a largura das estruturas (índice, etiquetas).  
Qual o número de bits (matriz de dados + etiquetas) de cada cache?
- 2) Considere um processador com relógio de 500MHz (2ns).  
Você dispõe de 96kbytes de memória com 1ns de tempo de acesso; 2048kbytes de memória com tempo de acesso de 10ns; 4Gbytes de RAM dinâmica com tempo de acesso de 60ns.  
Projete uma hierarquia de memória completa, que minimize o tempo médio de acesso à memória.  
Suponha que a memória necessária para implementar as etiquetas existe em abundância e com tempo de acesso adequado ao seu uso.

## Mais Exercícios

- 3) Considere os três barramentos do slide 26.  
Para blocos com 4, 8 e 16 palavras, calcule (i) a vazão de pico, (ii) a vazão sustentada (efetiva), para cada um dos três barramentos, considerando relógio de 1GHz (1ns). A unidade da vazão é [byte/s]
- 4) Suponha um programa que não causa faltas na L1-ins.  
A taxa de faltas na L1-D é 10%, e a taxa de faltas na L2 é 2%.  
O tempo de acesso à L1 é 1 ciclo, e o da L2 é 5 ciclos.  
A penalidade por falta na L2 é 100 ciclos.  
Qual o tempo médio de acesso à memória?  
Quais as taxas de falta locais em L1 e L2, e a taxa de faltas global?