

## 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

## Projeto de memórias cache

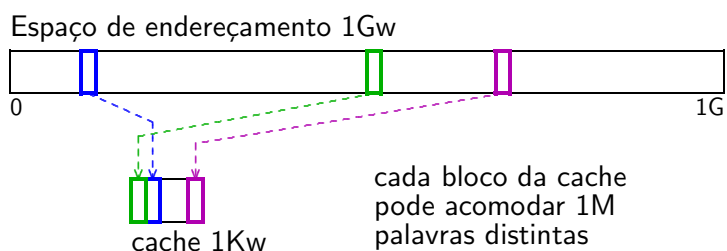
### Problema #1:

mapear espaço de endereçamento de 4Gbytes em 4Kbytes

necessita função que mapeia 1Gwords em 1Kwords

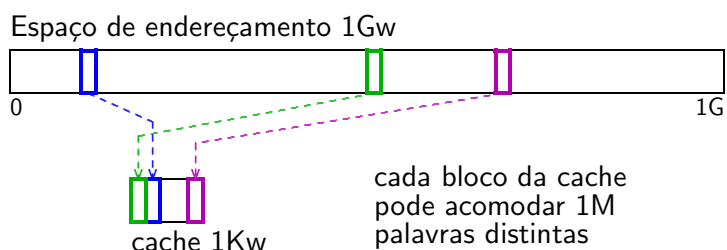
hashing

como implementa função em hardware, **eficientemente**?



## Projeto de memórias cache

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



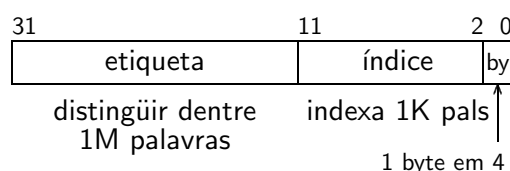
Mapeamento direto:

posição na cache

determinada por

10 bits do endereço

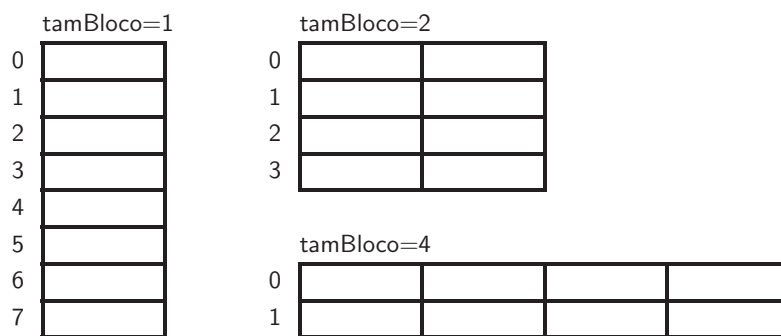
$pos = end\_WD \% 1024$



## Projeto de memórias cache

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

necessita acomodar mais de uma palavra em cada **bloco** da cache

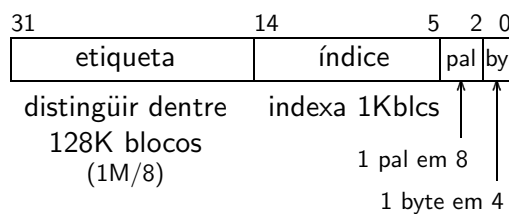


preencha as 3 caches com refs a palavras:

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

## Projeto de memórias cache

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



**Mapeamento direto:**

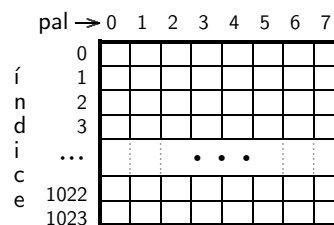
posição na cache

determinada por

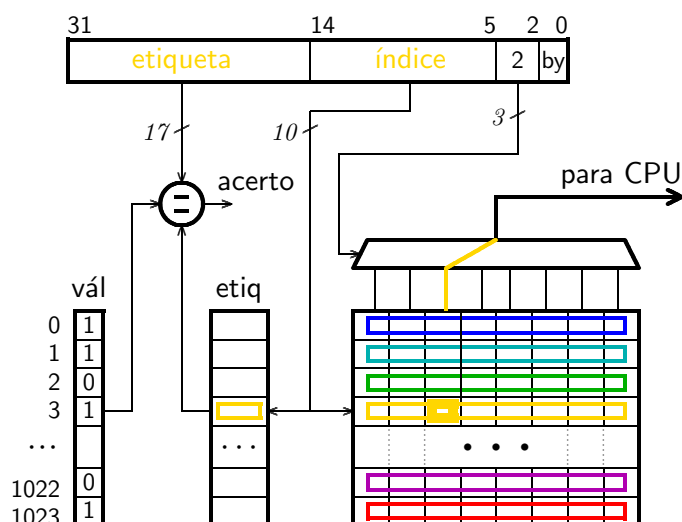
10 bits do endereço

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

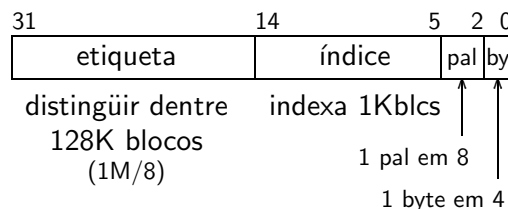
bloco acomoda 8 palavras de mesmo eti-q-índice



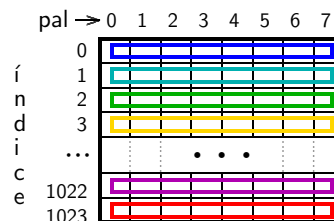
## Cache com Mapeamento Direto



## Mapeamento Direto e conflitos

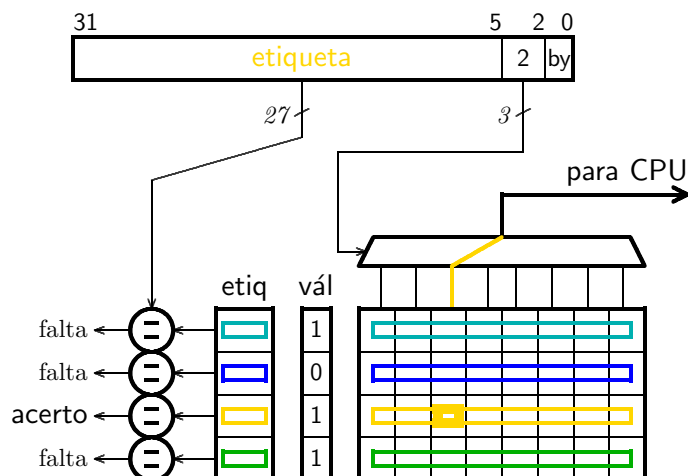


Mapeamento direto:  
cada bloco da cache  
pode armazenar palavras  
de uma “cor” mas  
memória contém 128K blocos  
de cada cor (nesta cache 1K,8pal/bl)



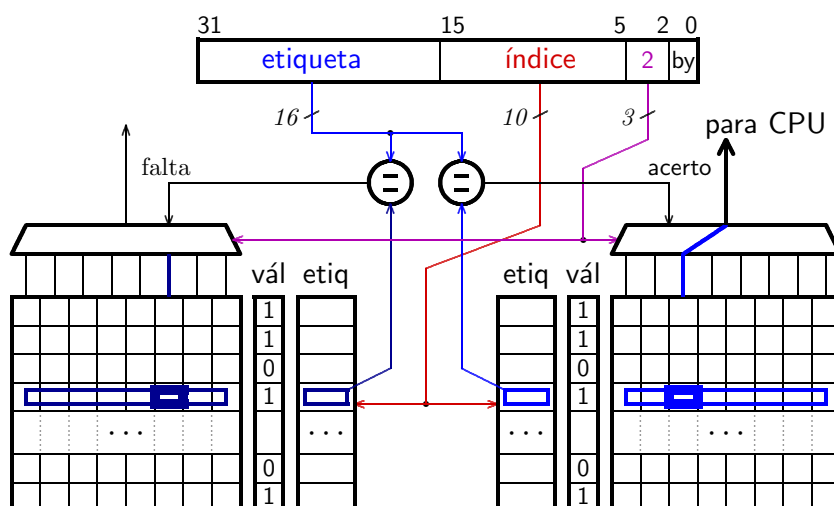
O que acontece se dois blocos da mesma cor são acessados  
em sequência? índice igual, etiquetas distintas

## Mapeamento Associativo



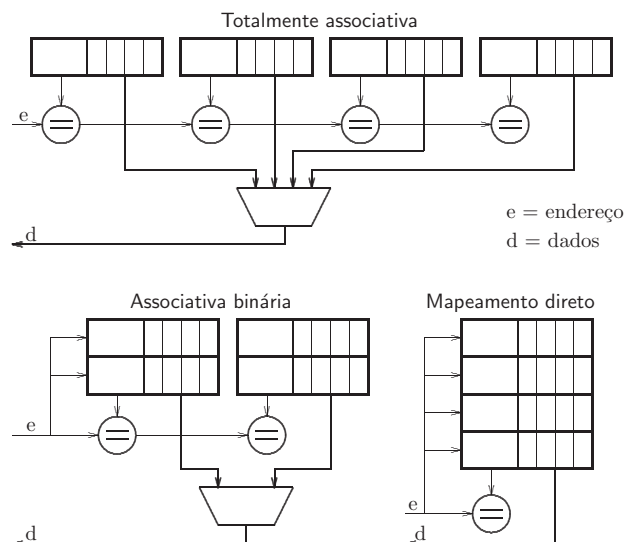
Mapeamento ideal e caro: um comparador em cada bloco da cache

## Mapeamento Híbrido: Associatividade por Conjuntos

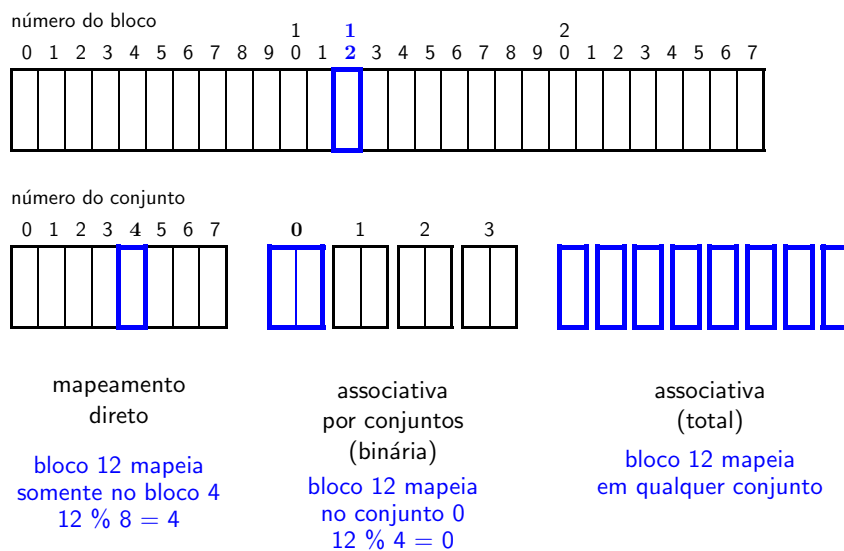


bloco de índice  $n$  pode estar em  $\forall$  elmt do conjunto  $N$   $|N|=2$

## Três Mapeamentos



## Alocação nos Blocos



## Alocação nos Blocos

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

Associativa - 4 conjuntos

0		99		200		33		
---	--	----	--	-----	--	----	--	--

alocação:  $12 \% 4 = ?$

Associativa binária - 2 conjuntos

0		99		33
1		200		

$12 \% 2 = 0$   
 $14 \% 2 = 0$   
 $13 \% 2 = 1$

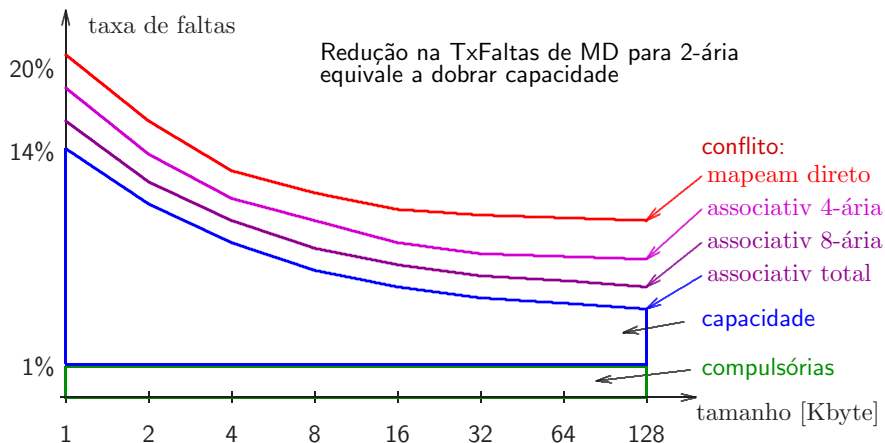
Mapeamento direto - 1 conjunto

0		99	$12 \% 4 = 0$
1		200	$13 \% 4 = 1$
2		33	$14 \% 4 = 2$
3			$xx \% 4 = 3$

Como ficam as etiquetas, em cada um dos 3 casos?

## Tipos de Faltas – falta por conflito

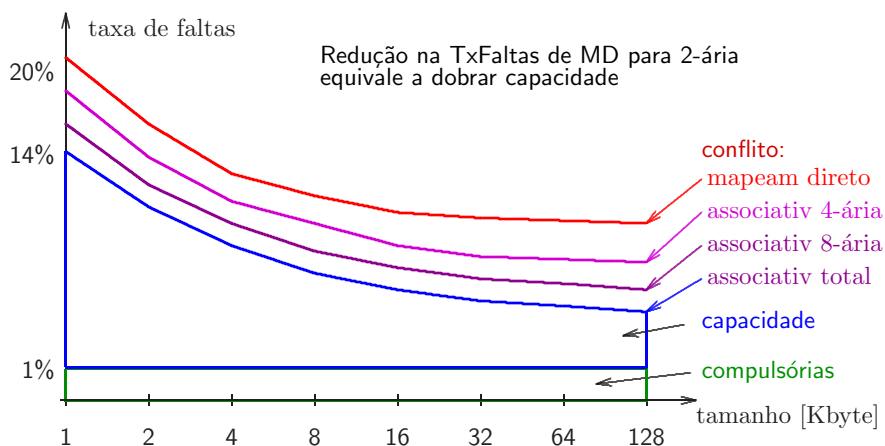
**problema:**  $N$  blocos da mesmo índice (côr) em uso simultâneo, mas cache só acomoda  $M < N$  blocos



**solução:** aumentar associatividade

## Tipos de Faltas – falta por capacidade

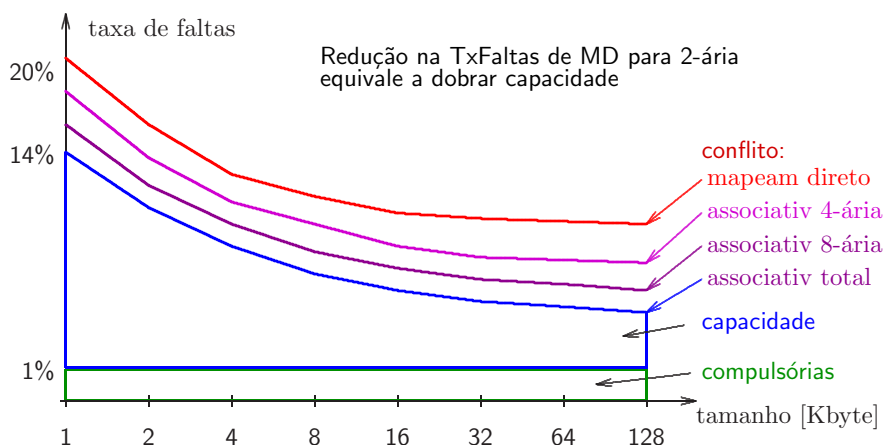
**problema:** cache não tem capacidade para acomodar todos blocos acessados por programa



**solução:** aumentar tamanho da cache

## Tipos de Faltas – falta compulsória

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



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

## Tipos de Faltas – 3 Cs

- 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)

## 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?**

**escolher aleatoriamente?**

Parece ser o melhor, mas  
pode ser difícil de implementar

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**

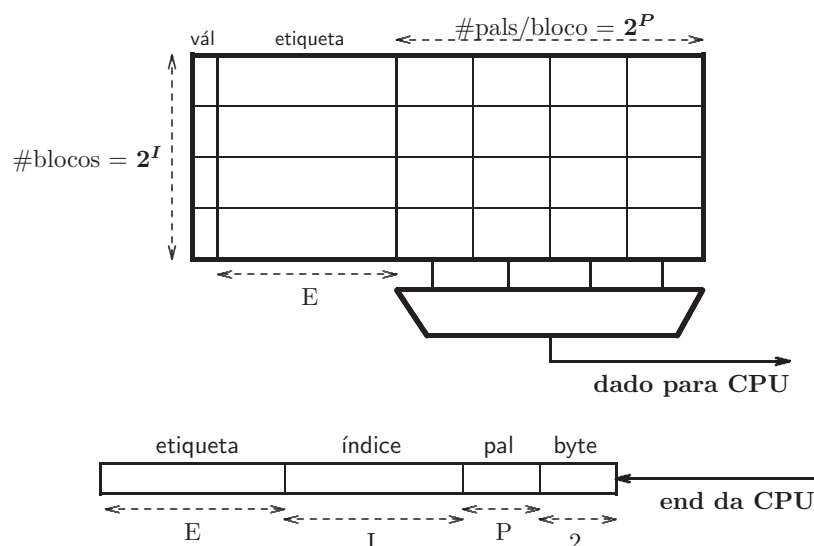
Implementação:

LRU 2-way: bit aponta para LRU

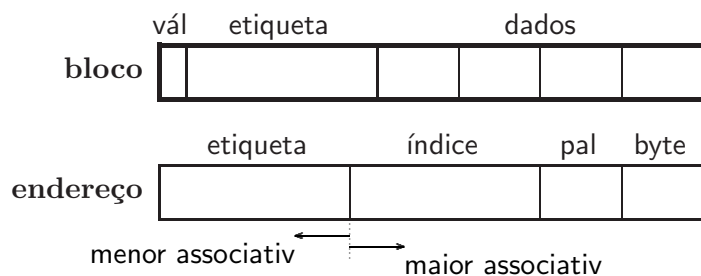
LRU 4-way: 3 bits, hierárquico

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

## Organização – Detalhes



## Organização – Detalhes



capacidade = num\_blocos  $\times$  tam\_bloco  $\times$  tam\_palavra  $\times$  associativ

$|\text{índice}| = \log_2 \text{ num\_blocos}$  altura da matriz

$|\text{pal}| = \log_2 \text{ tam\_bloco}$  # palavras/dados, largura

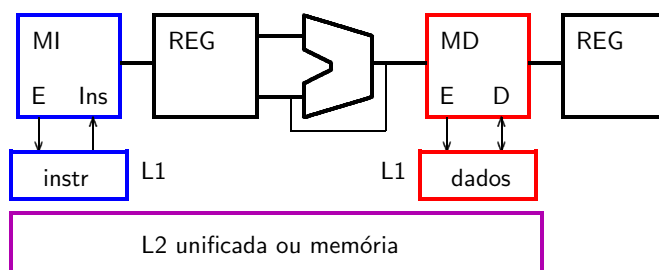
$|\text{byte}| = \log_2 \text{ tam\_palavra}$

$|\text{etiqueta}| = \text{tam\_ender} - ( (|\text{índice}|/\text{associativ}) + |\text{pal}| + |\text{byte}| )$

## Caches e Pipelines

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

Se muito longo  $\rightarrow$  mais de um ciclo para acessar a cache nos acertos



Há duas portas para memória para eliminar risco estrutural;

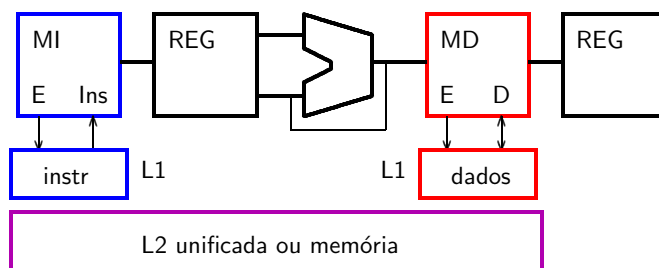
2 caches separadas, uma para instruções, uma para dados porque

localidade no acesso aos dados é diferente do acesso a instruções

## Caches e Pipelines

**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

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**

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

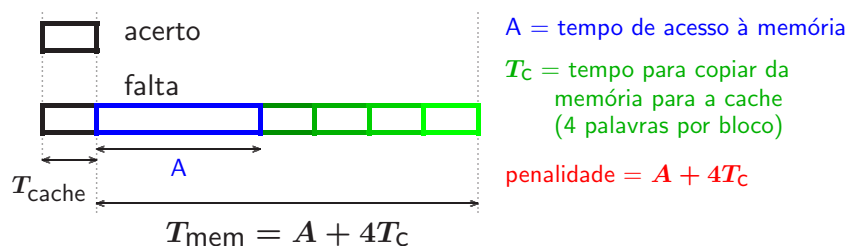
objetivo de projeto: **minimizar TMAM para dados e instruções**

**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 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 ↗

$$\text{capacidade} = \text{num\_blocos} \times \text{tam\_bloco} \times \text{tam\_palavra}$$

- \* aumentar associatividade

tempo de acesso ↗

$$\text{capacidade} = (\text{n\_blocos} \times \text{associativ}) \times \text{t\_bloco} \times \text{t\_palavra}$$

- \* busca antecipada por hardware ou por software

cuidado para não desalojar blocos úteis

poluição

traz blocos e armazena em buffer específico

$$\text{capacidade MD} = \text{num\_blocos} \times \text{tam\_bloco} \times \text{tam\_palavra}$$

$$[\text{byte}] = [\text{bloco}] \times [\text{pal/bloco}] \times [\text{byte/pal}]$$

## 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);

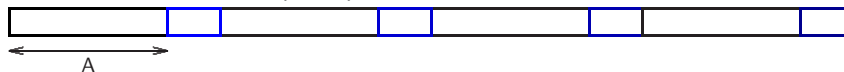
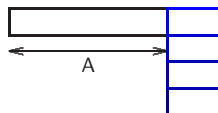
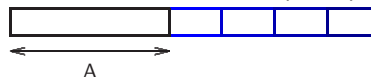
cache de vítimas: cache associativa

cache organizada em bancos;

pequena (4-8 blcs) mantém blocos

barramento CPU-cache com transações;

expurgados da L1

barramento estreito  $4(A+1)$ barramento largo  $(A+1)$ memória entrelaçada  $(A+4)$ 

## Minimizar Penalidade por Falta (ii)

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

Minimizar tempo de acerto

no nível inferior (L2 ou memória)

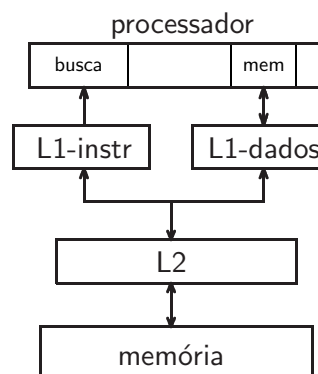
$$Tx_{\text{FaltasLocal}} L_n = \frac{\text{faltas-}L_n}{\text{refs-}L_n}$$

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

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)

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

- **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** → p3 → p0 → p1
  - \* localidade espacial indica que próximo acesso será na palavra seguinte, que pode não ter chegado ainda

## Resumo

- Hashing para mapear espaço grande em armazenador pequeno se largura aumenta ( $\text{tamBloco} > 1$ ) então altura diminui
- Três formas de hashing:
  - \* mapeamento direto → posição =  $\text{ender \% num\_blocos}$
  - \* totalmente associativo → posição =  $\forall$  no conjunto (índice)
  - \* associativo por conjuntos → posição =  $\text{ender \% num\_conj}$
- Três categorias de faltas:
  - \* por capacidade → cache maior
  - \* por conflitos → associatividade
  - \* compulsórias → 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, e as associatividades indicadas abaixo, com blocos de 32 bytes. (a) mapeam direto, (b) associativ binária, (c) associativ 8-ária. Indique claramente a largura das estruturas (índice, etiquetas). Qual o número de bits total (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 (L1,L2) e global?