



Hierarquia de memória

Memórias Cache

Cap. 5: “Large and Fast: Exploiting Memory Hierarchy”

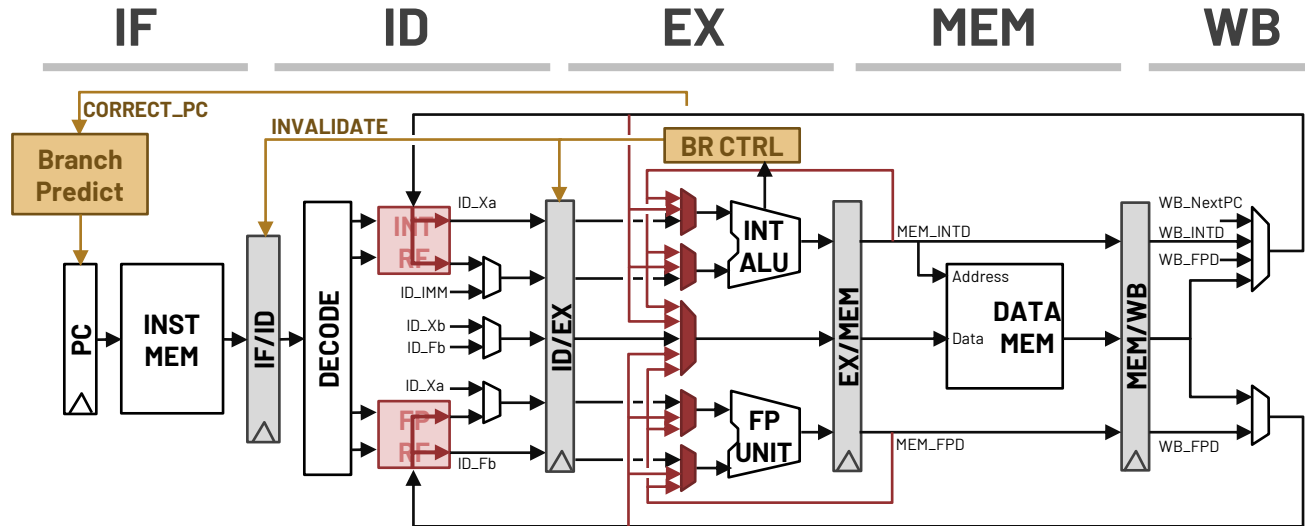


Desenho do processador

Revisão do capítulo anterior

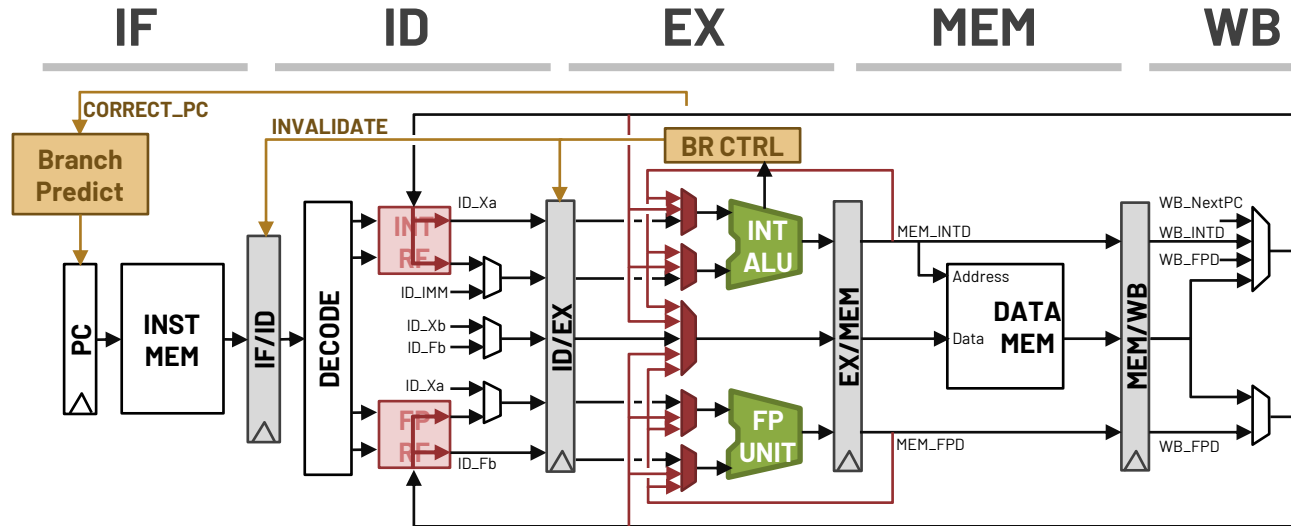
Arquitetura do processador (core)

Funcionamento em pipeline com caminhos de *forwarding* e predição de salto



Arquitetura do processador (core)

Funcionamento em pipeline com caminhos de *forwarding* e predição de salto

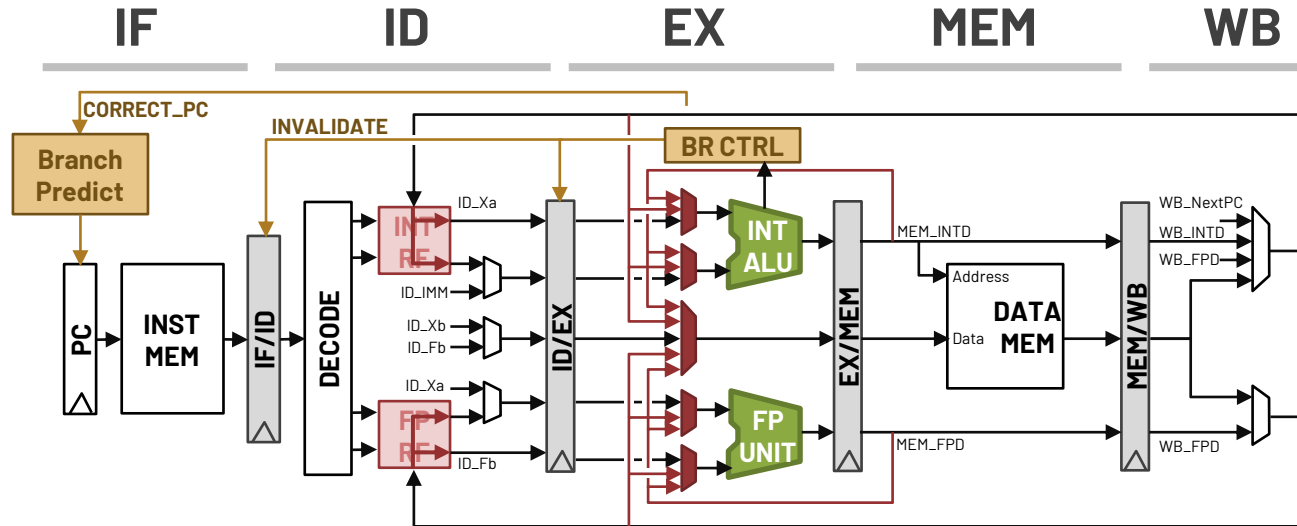


As unidades funcionais podem ser implementadas com pipeline, e cada uma, com um número de estágios diferente

Unidade funcional	#ciclos
Unidade lógica	1
Unidade de deslocamento	1
Somador/subtrator (inteiros)	1
Multiplicador	3
Divisor (inteiro/FP)	20
Somador/subtrator (FP)	3
Multiplicador (FP)	4

Arquitetura do processador (core)

Funcionamento em pipeline com caminhos de *forwarding* e predição de salto



As unidades funcionais podem ser implementadas com pipeline, e cada uma, com um número de estágios diferente

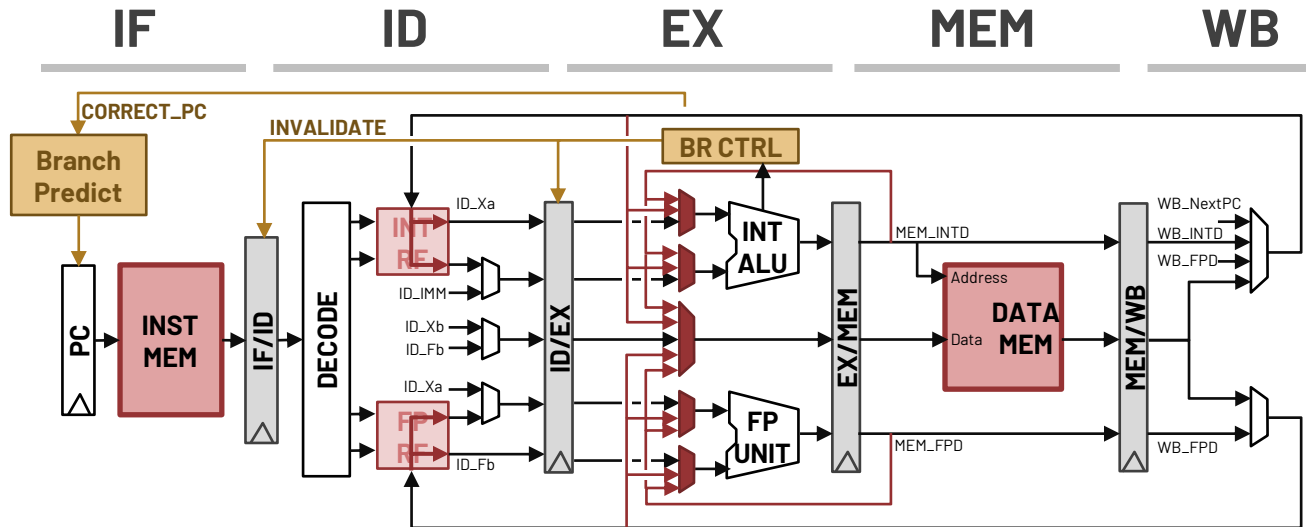
Unidade funcional	#ciclos
Unidade lógica	1
Unidade de deslocamento	1
Somador/subtrator (inteiros)	1
Multiplicador	3
Divisor (inteiro/FP)	20
Somador/subtrator (FP)	3
Multiplicador (FP)	4

Nota: nestes casos as instruções que não acedem à memória tipicamente fazem bypass ao estágio de MEM

Arquitetura do processador (core)

Funcionamento em pipeline com caminhos de *forwarding* e predição de salto

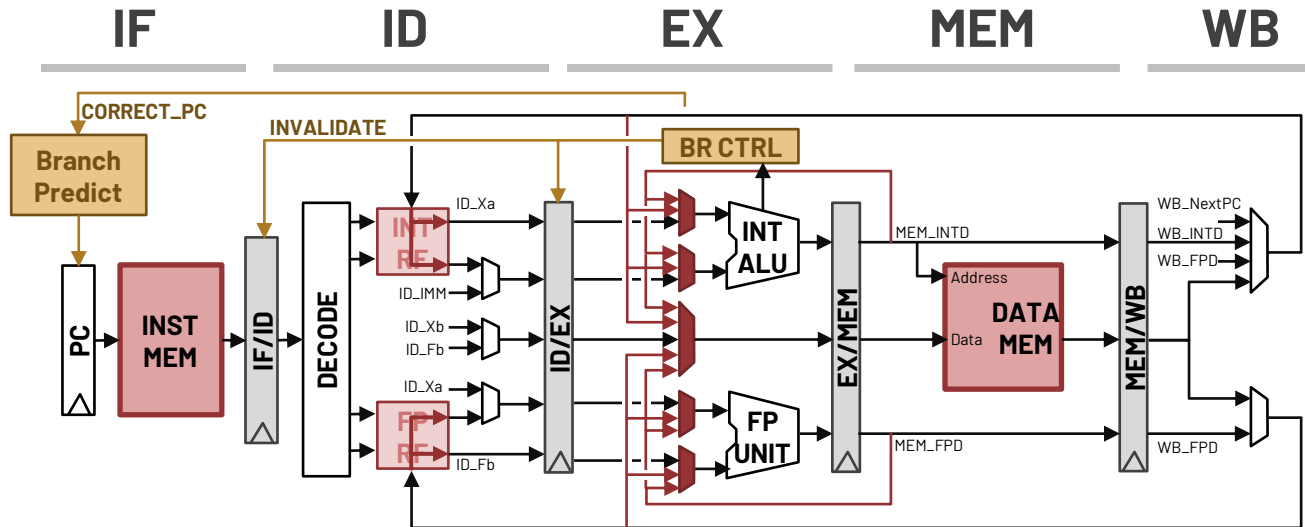
- Duas memórias?
- Leitura assíncrona com baixa latência?



Arquitetura do processador (core)

Funcionamento em pipeline com caminhos de *forwarding* e predição de salto

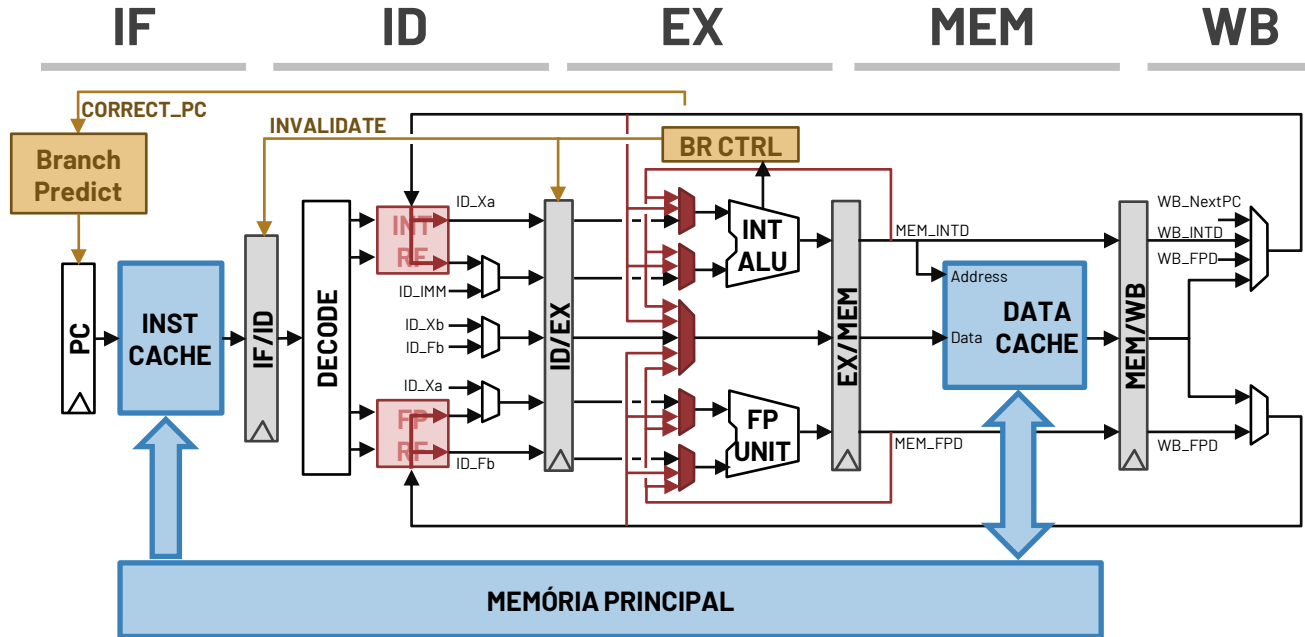
- Duas memórias? Geralmente, existe apenas uma memória principal (memória RAM)!
- Leitura assíncrona com baixa latência? Latência demasiado grande (~100 ciclos)!



Arquitetura do processador (core)

Funcionamento em pipeline com caminhos de *forwarding* e predição de salto

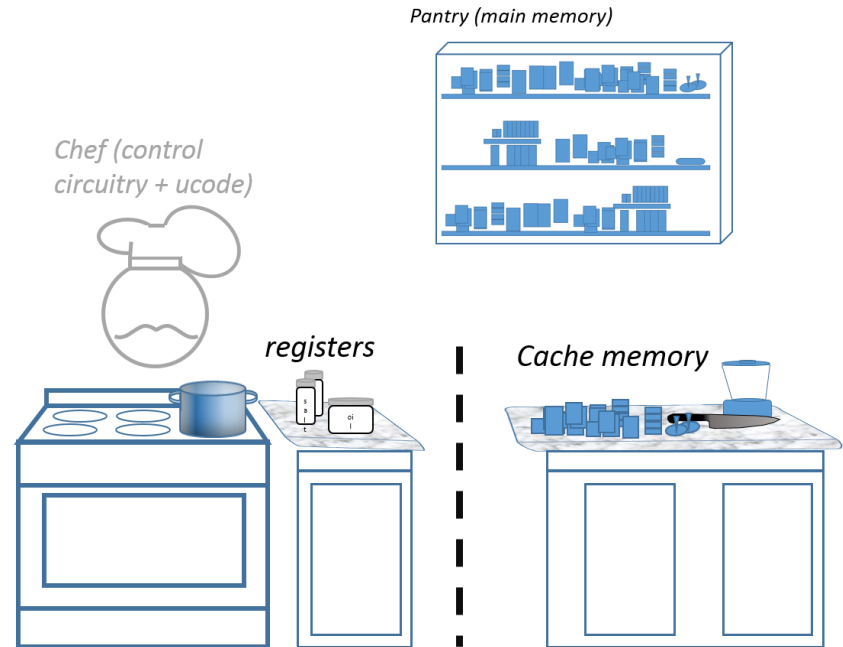
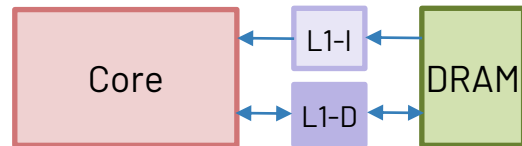
- Duas caches (pequenas → baixa latência) que ligam à mesma memória principal
- Cada cache contém apenas um subconjunto das instruções ou dados



Arquitetura do processador

Hierarquia de memória

- Duas caches (pequenas → baixa latência) que ligam à mesma memória principal
- Cada cache contém apenas um subconjunto das instruções ou dados

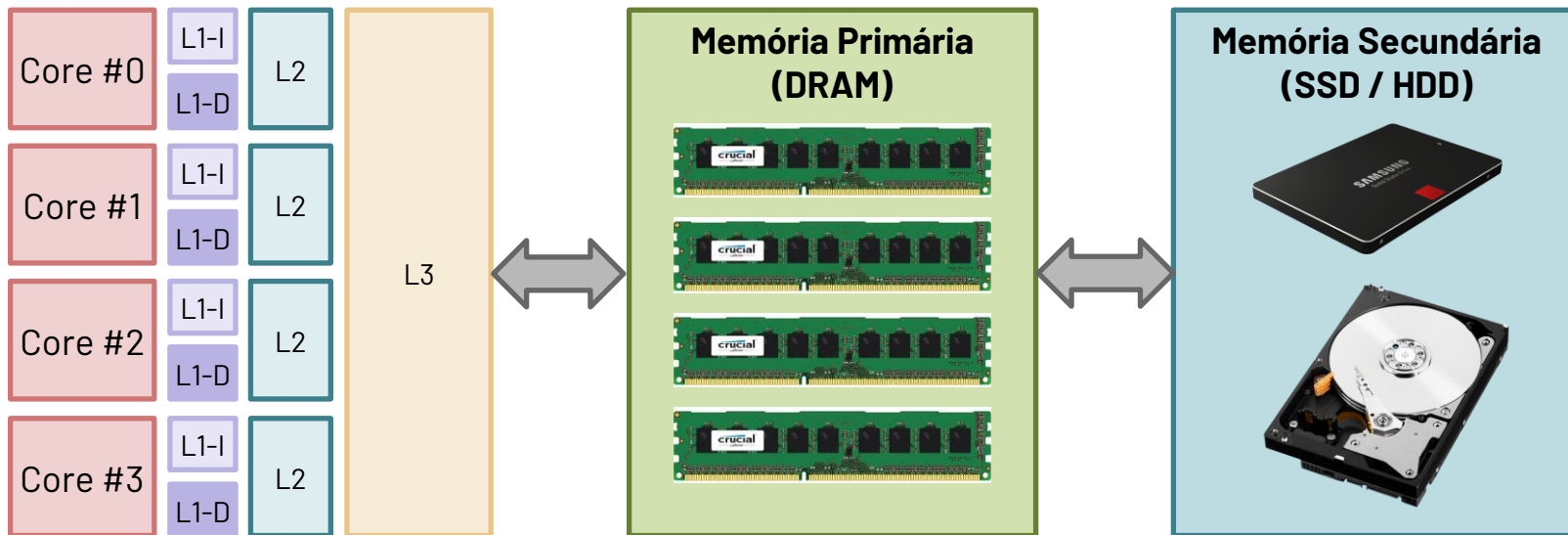


- Duas caches (pequenas → baixa latência) que ligam à mesma memória principal
 - Embora as caches sejam de dimensão reduzida (valores típicos 32kB a 64kB), a latência é ainda assim demasiado grande para suportar frequências de funcionamento muito altas. Assim, o acesso à cache é tipicamente partido em 3-4 ciclos (i.e., para o caso de uma latência de 3 ciclos, o pipeline tem mais estágios correspondentes a IF1,IF2,IF3 e MEM1,MEM2,MEM3).
- Cada cache contém apenas um subconjunto das instruções ou dados
 - Se os dados que queremos consultar estiverem em cache, não é preciso pagar o “custo” de consultar o nível superior, ou a memória RAM

Arquitetura do processador

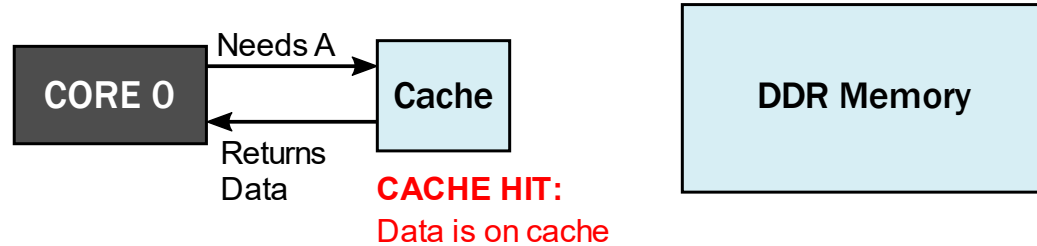
Hierarquia de memória

- Como as cache L1 têm uma dimensão reduzida, geralmente os sistemas contêm uma hierarquia com 2 a 3 níveis de cache (ex: L1, L2, L3). A cache L1 é a de menor dimensão e menor latência, portanto mais próxima do processador

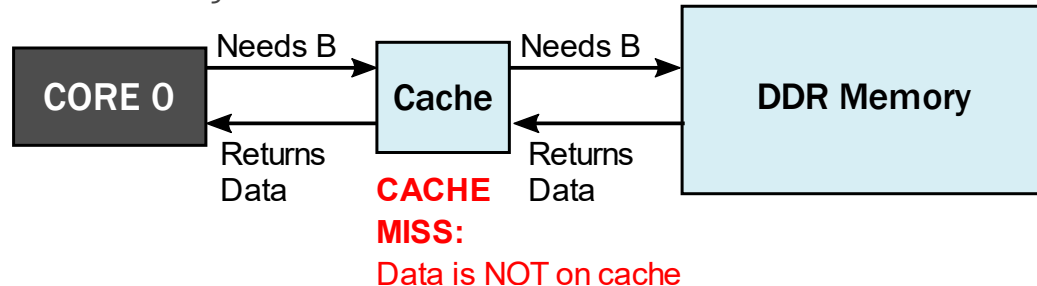


Princípio de funcionamento

- Quando o processador precisa de um dado, faz o pedido à cache
 - Se os dados estiverem na cache, devolve-os imediatamente



- Se os dados NÃO estiverem na cache, faz o pedido à memória
 - Quando o nível seguinte devolver os dados,





Princípio da localidade

- **Localidade temporal:**

Se acedermos (leitura ou escrita) a um endereço A, é provável que num futuro próximo necessitemos de aceder novamente ao mesmo endereço

- E.g., sequência de instruções em ciclos e acessos frequentes a estruturas de dados

- **Localidade espacial:**

Se acedermos (leitura ou escrita) a um endereço A, é provável que num futuro próximo necessitemos de aceder aos endereços adjacentes a A, i.e., A-1, A+1, A+2, ...

- E.g., execução de código (sequência de instruções, na ausência de saltos), ou acessos a elementos num vetor ou outra estrutura de dados

- **Regra 90/10**

Um programa passa 90% do tempo a executar 10% das instruções.

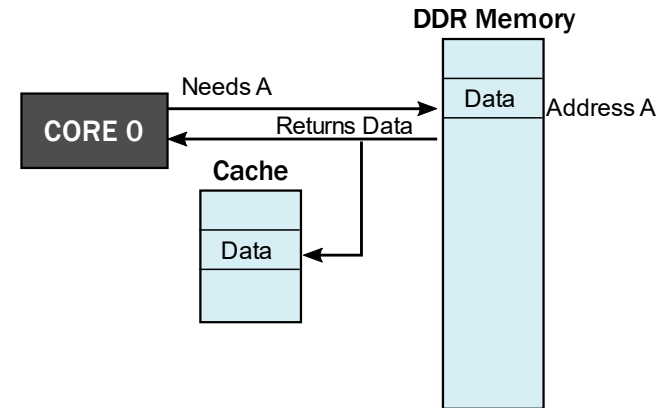
- **Localidade temporal:**

Se acedermos (leitura ou escrita) a um endereço A, é provável que num futuro próximo necessitemos de aceder novamente ao mesmo endereço

- E.g., sequência de instruções em ciclos e acessos frequentes a estruturas de dados



Quando acedemos ao dado/instrução no endereço A, já que temos de pagar o custo de o ir buscar à memória RAM, mais vale armazená-lo em memória cache.



FIRST ACCESS:

Data access time: ~100 cycles

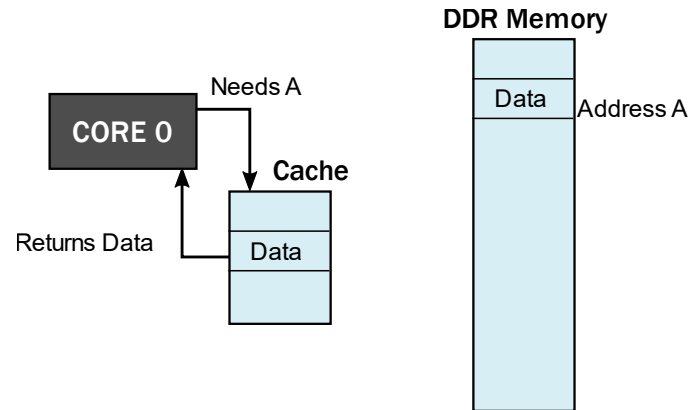
- **Localidade temporal:**

Se acedermos (leitura ou escrita) a um endereço A, é provável que num futuro próximo necessitemos de aceder novamente ao mesmo endereço

- E.g., sequência de instruções em ciclos e acessos frequentes a estruturas de dados



Quando acedemos ao dado/instrução no endereço A, já que temos de pagar o custo de o ir buscar à memória RAM, mais vale armazená-lo em memória cache.



FOLLOWING ACCESS:

Data access time: ~3 cycles

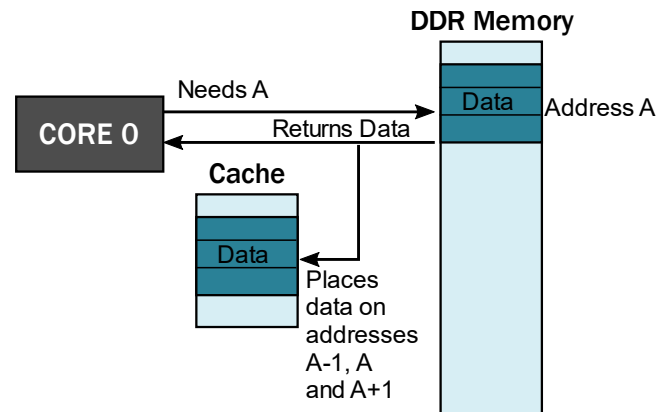
■ Localidade espacial:

Se acedermos (leitura ou escrita) a um endereço A, é provável que num futuro próximo necessitemos de aceder aos endereços adjacentes a A, i.e., A-1, A+1, A+2, ...

- E.g., execução de código (sequência de instruções, na ausência de saltos), ou acessos a elementos num vetor ou outra estrutura de dados



Quando acedemos ao dado/instrução no endereço A, já que temos de pagar o custo de o ir buscar à memória RAM, mais vale trazer também alguns dados/instruções contíguas.



FIRST ACCESS:

Data access time: ~100 cycles

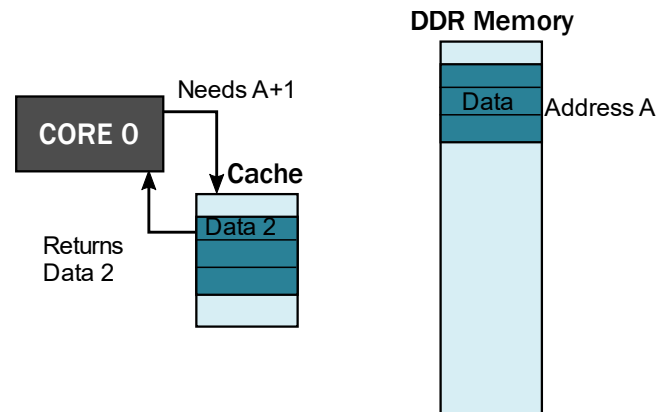
■ Localidade espacial:

Se acedermos (leitura ou escrita) a um endereço A, é provável que num futuro próximo necessitemos de aceder aos endereços adjacentes a A, i.e., A-1, A+1, A+2, ...

- E.g., execução de código (sequência de instruções, na ausência de saltos), ou acessos a elementos num vetor ou outra estrutura de dados



Quando acedemos ao dado/instrução no endereço A, já que temos de pagar o custo de o ir buscar à memória RAM, mais vale trazer também alguns dados/instruções contíguas.



FOLLOWING ACCESS:

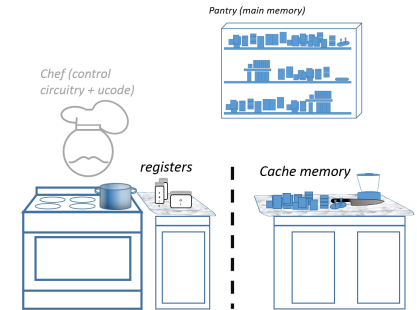
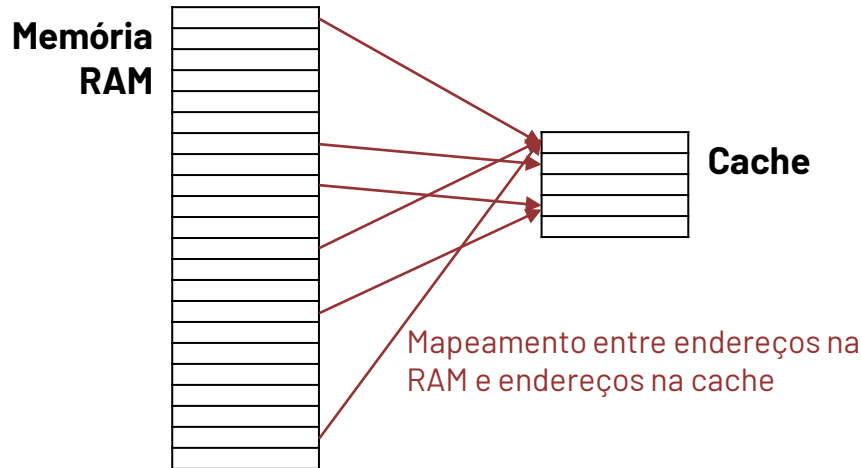
Data access time: ~3 cycles

A black and white photograph of a server room. On the left, there are several rows of server racks with glass doors, receding into the distance. The floor is made of large, light-colored tiles. The ceiling has recessed lighting. A large blue rectangular box is overlaid on the right side of the image, containing the title text in white.

Organização das memórias cache

Princípio de funcionamento da cache

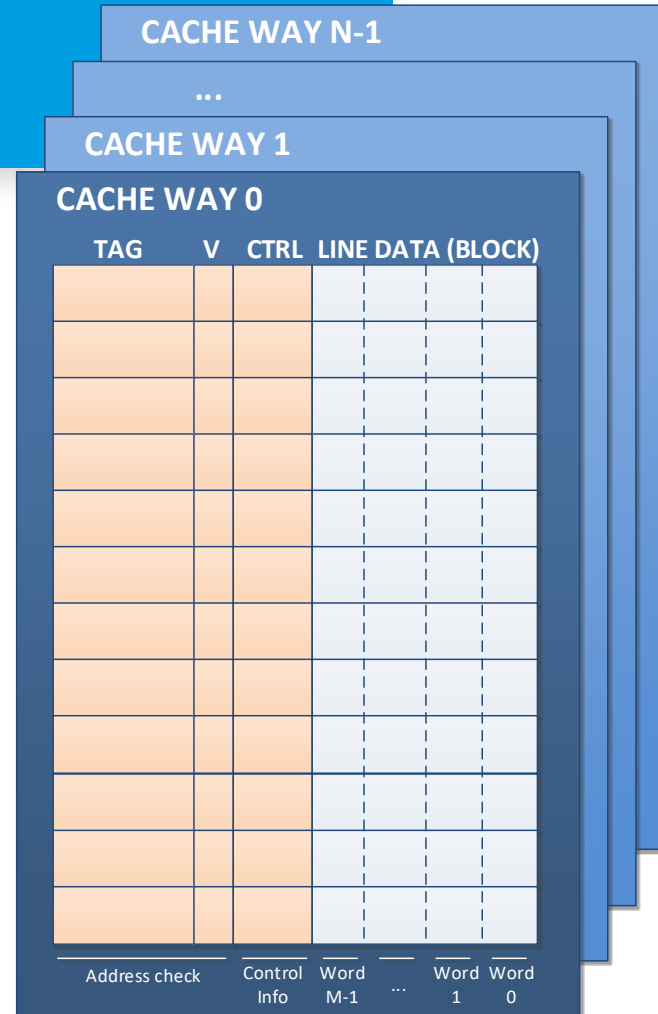
- Na prática, a cache mapeia os endereços da memória RAM (ex: 16GB) num conjunto limitado de entradas (ex: 1024 entradas)
- Naturalmente, este mapeamento dá origem a que múltiplas entradas da memória RAM correspondam à mesma entrada na cache.



- Na prática, a cache mapeia os endereços da memória RAM (ex: 16GB) num conjunto limitado de entradas (ex: 1024 entradas)
- Naturalmente, este mapeamento dá origem a que múltiplas entradas da memória RAM correspondam à mesma entrada na cache.
- Assim, para além de ser necessário guardar na cache os **dados**, é ainda preciso guardar uma **etiqueta** (TAG) que permite identificar se uma entrada na cache corresponde ao endereço pedido (A), ou a outro endereço (ex: A2, A3, A4, ...)
- Para explorar a localidade espacial, os blocos de dados na cache têm uma dimensão igual ou superior à das palavras do processador (i.e., $\geq 4B$ para RV32, $\geq 8B$ para RV64)

Organização das memórias cache

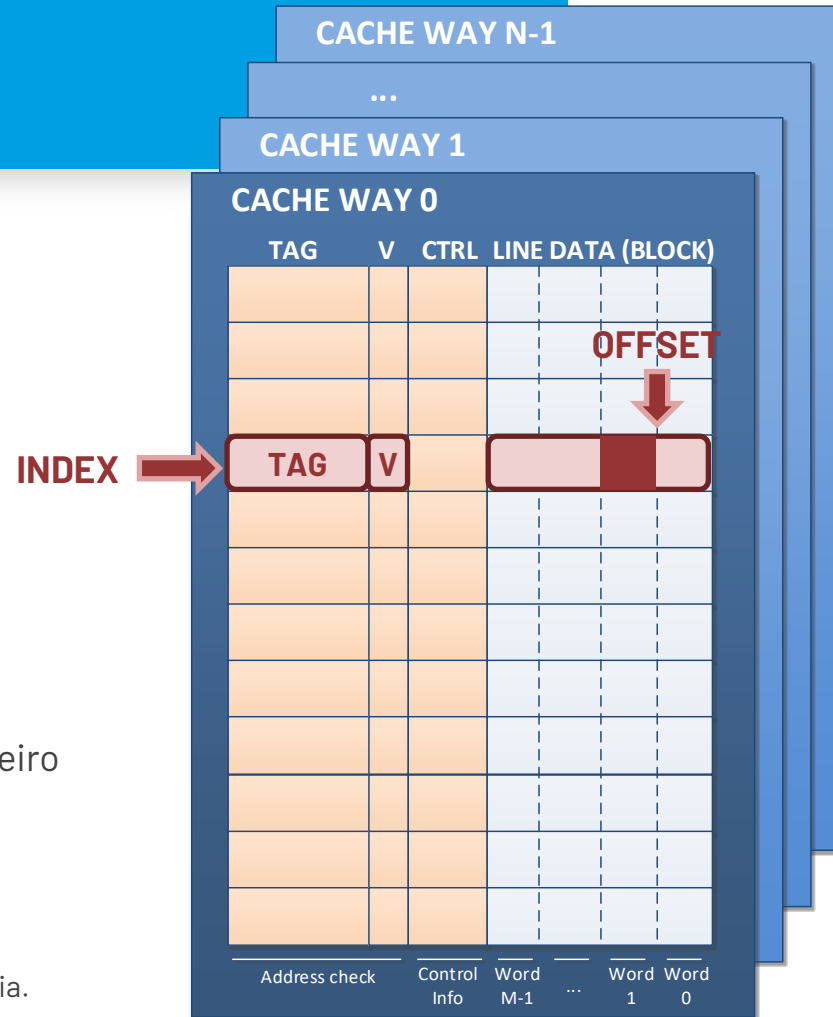
- Casos particulares:
 - **Cache de Mapeamento Direto:**
Apenas uma via, mas com múltiplas entradas (linhas)
 - **Cache associativa**
N Vias, cada uma com múltiplas linhas.
A organização de cada uma das N vias é idêntica
 - **Cache completamente associativa**
Múltiplas vias, mas cada uma contendo apenas uma linha



Organização das memórias cache

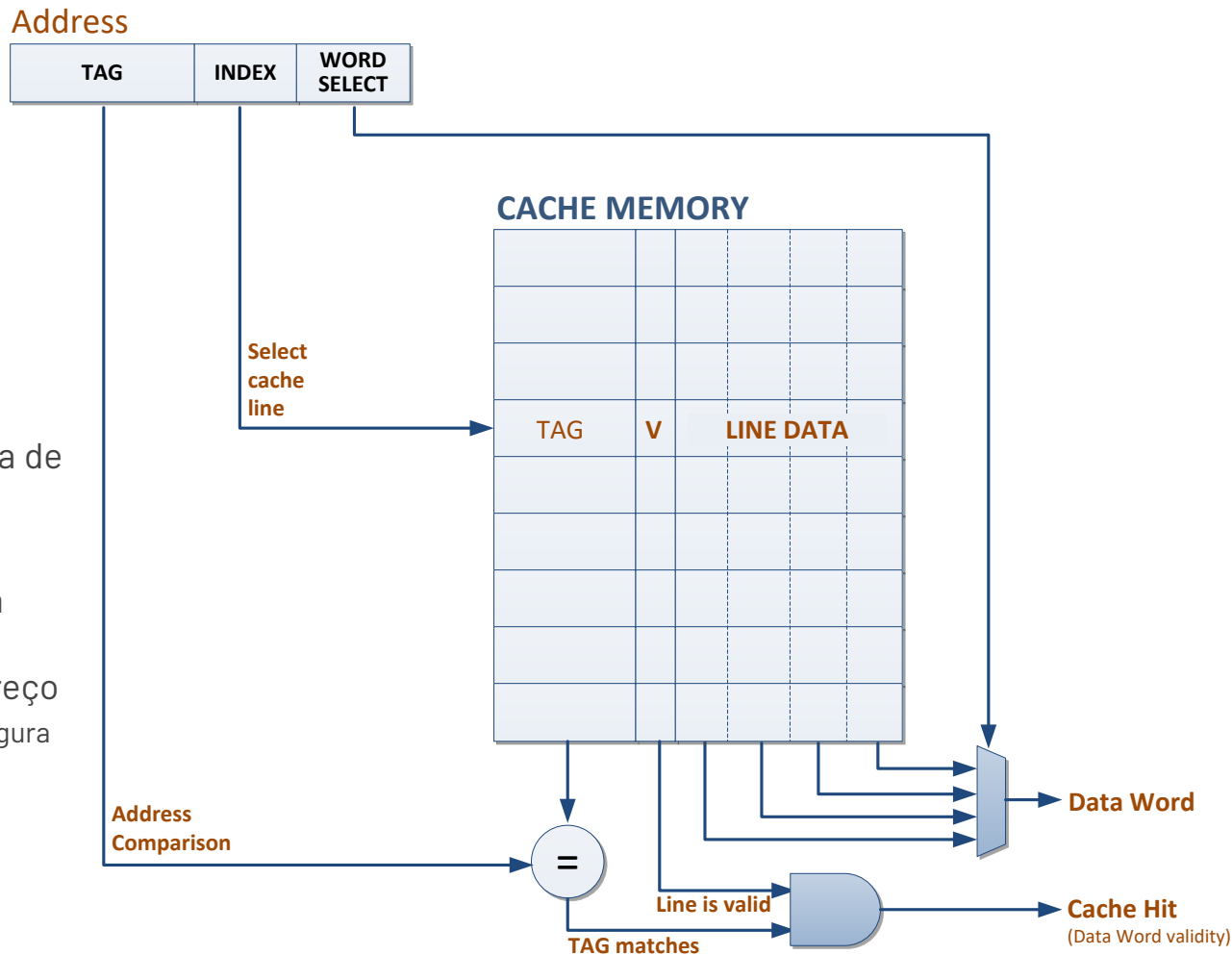
- Para realizar a tradução entre o endereço de memória (A) e localização na cache, é necessário:
 1. Escolher a entrada na(s) tabela(s), i.e., **INDEX**
O índice (ou INDEX) funciona como “endereço” na cache
 2. Verificar se a entrada é válida, i.e., se $V=1$
 3. Verificar se a entrada corresponde ao endereço pedido, i.e., confirmar o valor da **TAG**
 4. (se 2 e 3 forem verdadeiro → **HIT**)
ler a palavra pedida a partir da cache, sendo o primeiro byte indicado pelo valor de **OFFSET**

Nota: em caches associativas (i.e., quando o número de vias $N > 1$), é necessário verificar se os pontos 2 e 3 são verdadeiros em alguma das N vias, e devolver a palavra (a partir do valor de OFFSET) dessa via. Por definição, o valor do endereço de memória A só pode estar numa única via.



Mapeamento dos dados na cache

- Quando o processador precisa de um dado (ou instrução), faz o pedido à cache L1-D (ou L1-I)
- De acordo com a estrutura da cache L1-D (ou L1-I), esta decompõe a palavra de endereço em TAG, INDEX, OFFSET (na figura WORD SELECT)



Mapeamento dos dados na cache

- Quando o processador precisa de um dado (ou instrução), faz o pedido à cache L1-D (ou L1-I)
- De acordo com a estrutura da cache L1-D (ou L1-I), esta decompõe a palavra de endereço em TAG, INDEX, OFFSET (na figura WORD SELECT)

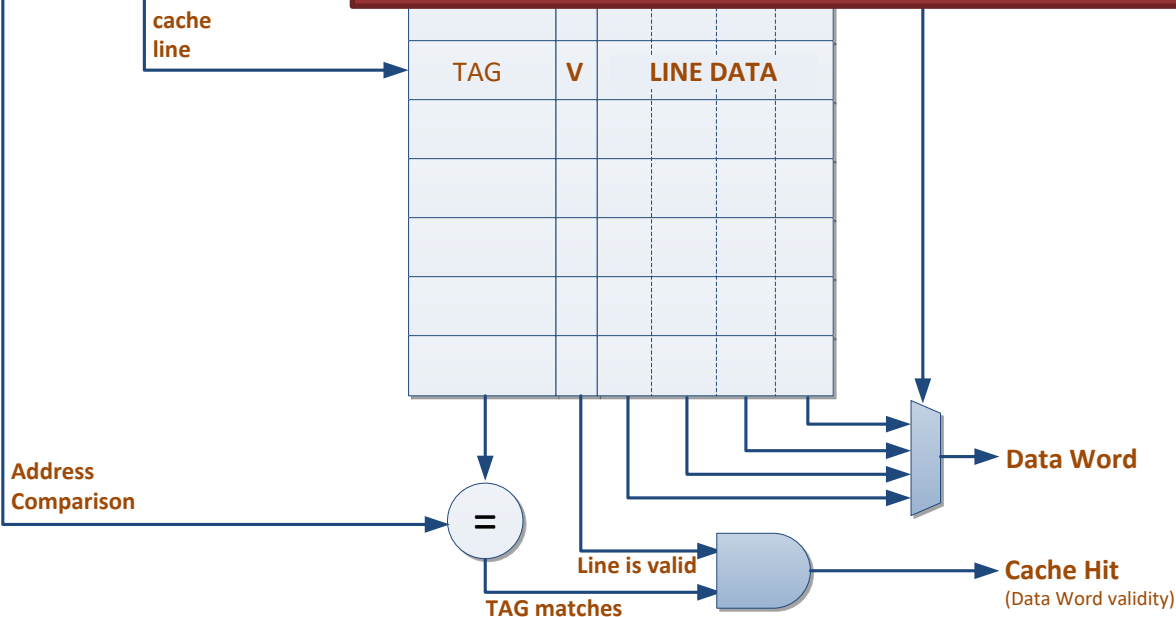
Address



Select cache line

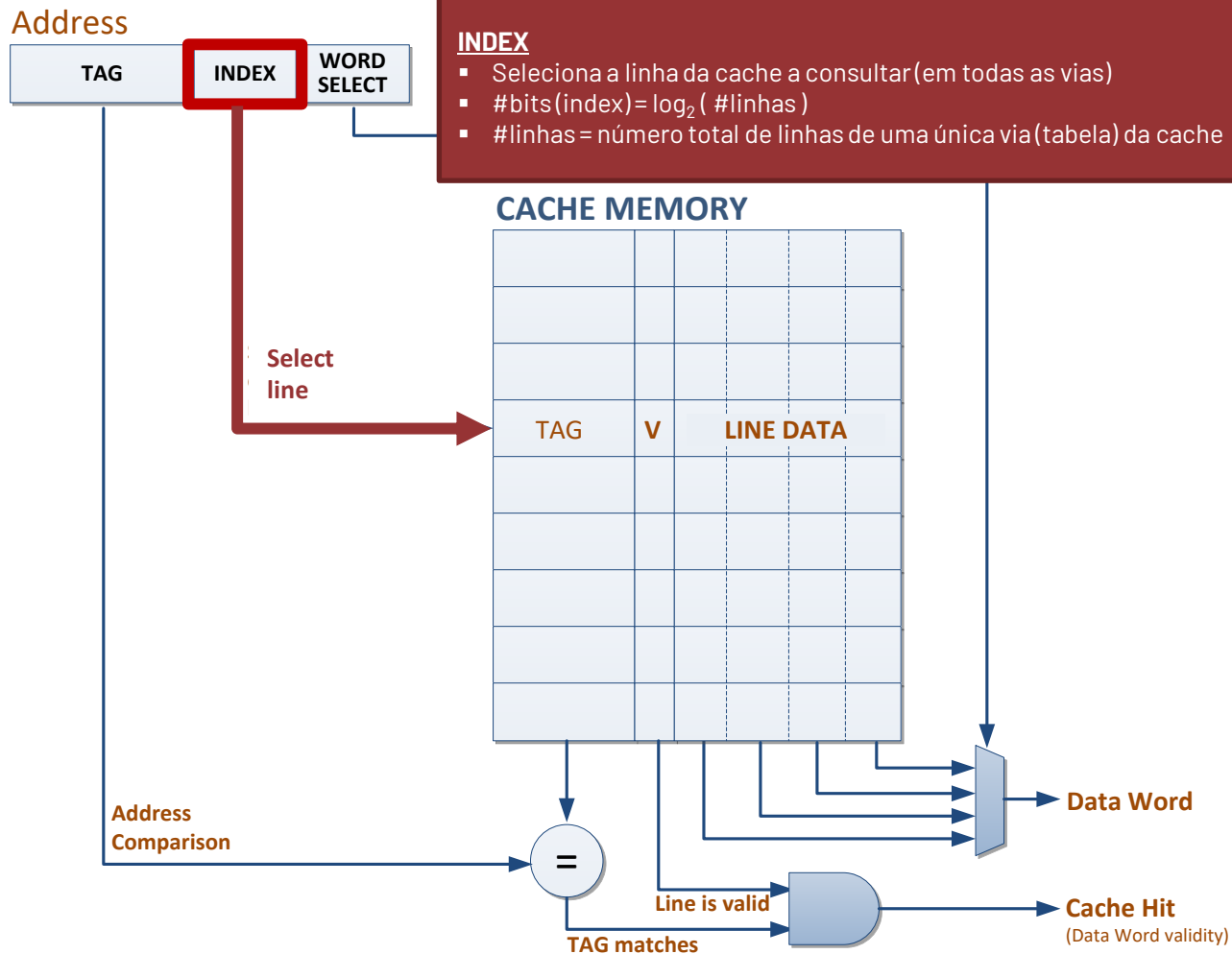
ADDRESS

- Por omissão, o campo de endereço geralmente tem a mesma dimensão que a largura do registo físico.
 - RV32 → registos de inteiros de 32 bits → endereços de 32 bits
 - RV64 → registos de inteiros de 64 bits → endereços de 64 bits
- Contudo, em alguns processadores a palavra de endereço é diferente da dimensão do registo, p. ex.:
 - Intel IA32 → registos de inteiros de 32 bits → endereços físicos de 48 bits
 - ARMv8 → registos de inteiros de 64 bits → endereços físicos de 52 bitsNestes casos o ISA especifica o valor exato do endereço físico.



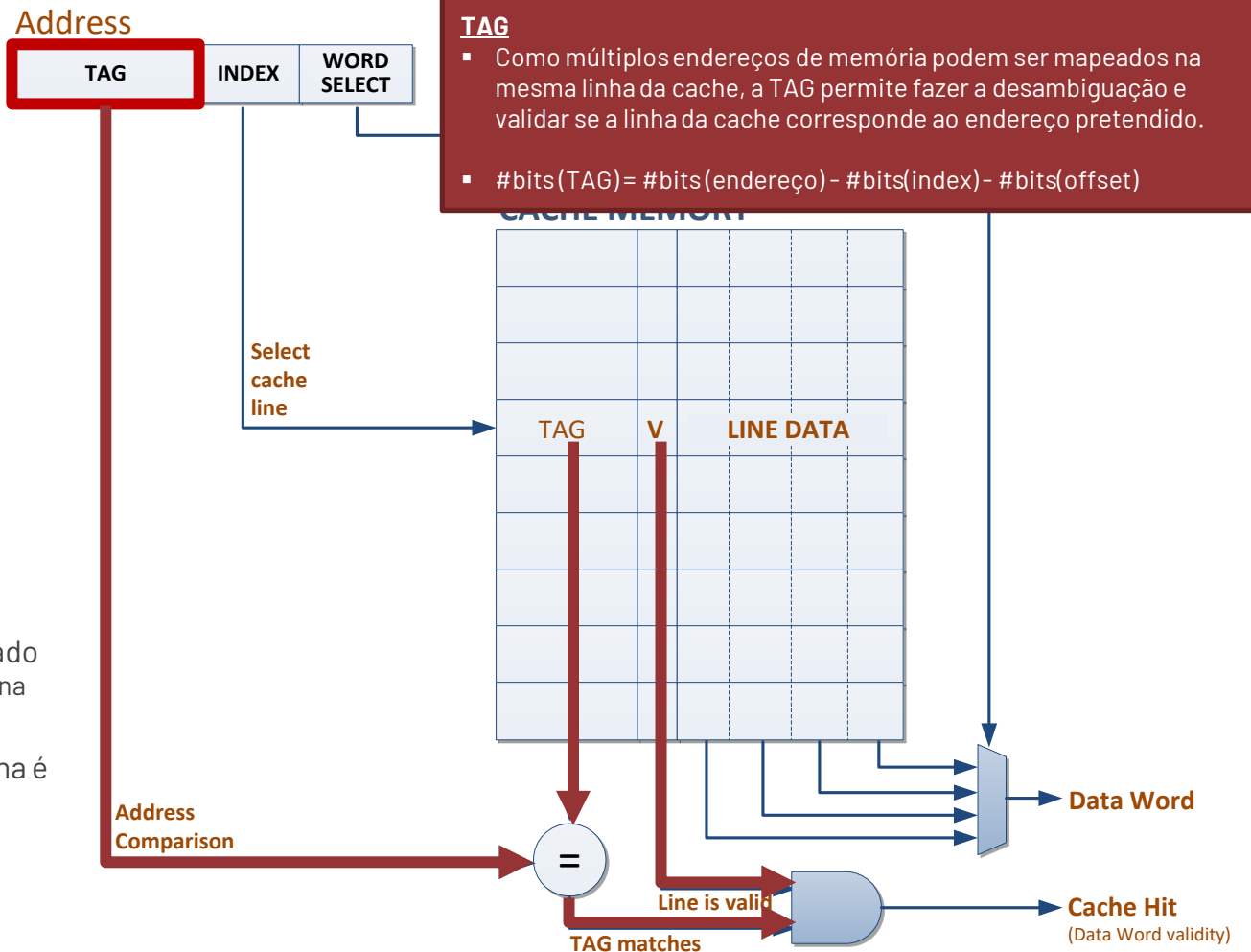
Mapeamento dos dados na cache

- O **INDEX** (índice) é usado para determinar a linha da cache que devemos consultar



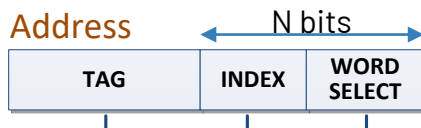
Mapeamento dos dados na cache

- O **INDEX** (índice) é usado para determinar a linha da cache que devemos consultar
- A **TAG** (etiqueta) valida se os dados/instruções na cache correspondem ao endereço indicado
 - Há múltiplos endereços mapeados na mesma linha
- O bit de validade (**V**) indica se a linha é válida

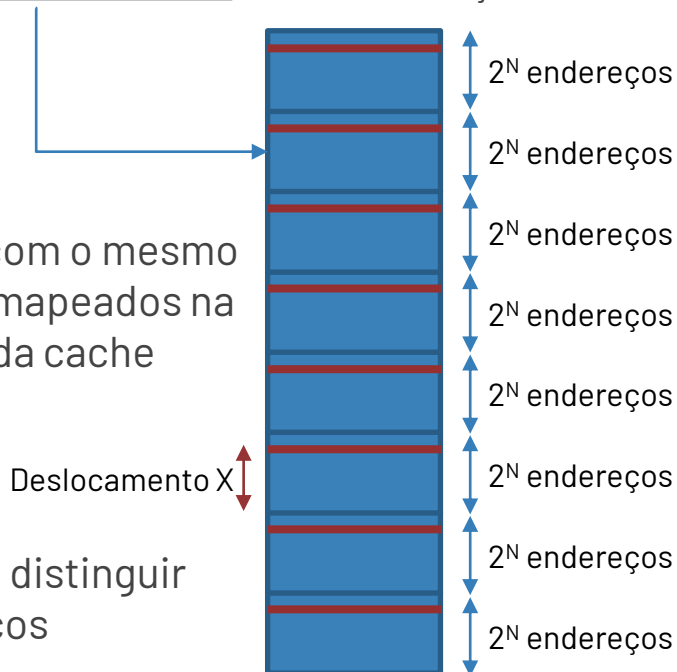


Mapeamento dos dados na cache

- O **INDEX** (índice) é usado para determinar a linha da cache que devemos consultar
- A **TAG** (etiqueta) valida se os dados/instruções na cache correspondem ao endereço indicado
 - Há múltiplos endereços mapeados na mesma linha
- O bit de validade (**V**) indica se a linha é válida



- O conjunto de N bits correspondentes ao conjunto dos campos INDEX+OFFSET definem um deslocamento dentro de um **espaço de memória** de 2^N endereços:



- Todos os endereços com o mesmo deslocamento X, são mapeados na mesma linha da linha da cache
- O campo TAG permite distinguir os diferentes endereços

Mapeamento dos dados na cache

- O **INDEX** (índice) é usado para determinar a linha da cache que devemos consultar
- A **TAG** (etiqueta) valida se os dados/instruções na cache correspondem ao endereço indicado
 - Há múltiplos endereços mapeados na mesma linha
- O bit de validade (**V**) indica se a linha é válida

Address



Select
cache
line

TAG

- Como múltiplos endereços de memória podem ser mapeados na mesma linha da cache, a TAG permite fazer a desambiguação e validar se a linha da cache corresponde ao endereço pretendido.
- $\#bits(TAG) = \#bits(endereço) - \#bits(index) - \#bits(offset)$

RESULTADO DO ACESSO À CACHE:

- Em caches associativas, é preciso verificar se alguma das vias deu HIT.
 - Por construção ou dá MISS em todas, ou dá HIT apenas numa das vias
- Existe um HIT se e só se:
 $LINE\ TAG = ADDRESS\ TAG \quad e \quad V=1$
- Se der MISS, temos de ir buscar *TODA* a linha à memória (ou nível de cache seguinte), e coloca-la na cache (numa das vias)

Address
Comparison



Line is valid

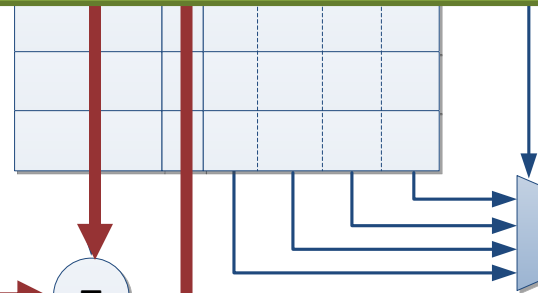
TAG matches



Cache Hit

(Data Word validity)

Data Word



Mapeamento dos dados na cache

- Quando o processador precisa de um dado (ou instrução), faz o pedido à cache L1-D (ou L1-I)
- De acordo com a estrutura da cache L1-D (ou L1-I), esta decompõe a palavra de endereço em TAG, INDEX, OFFSET (na figura WORD SELECT)

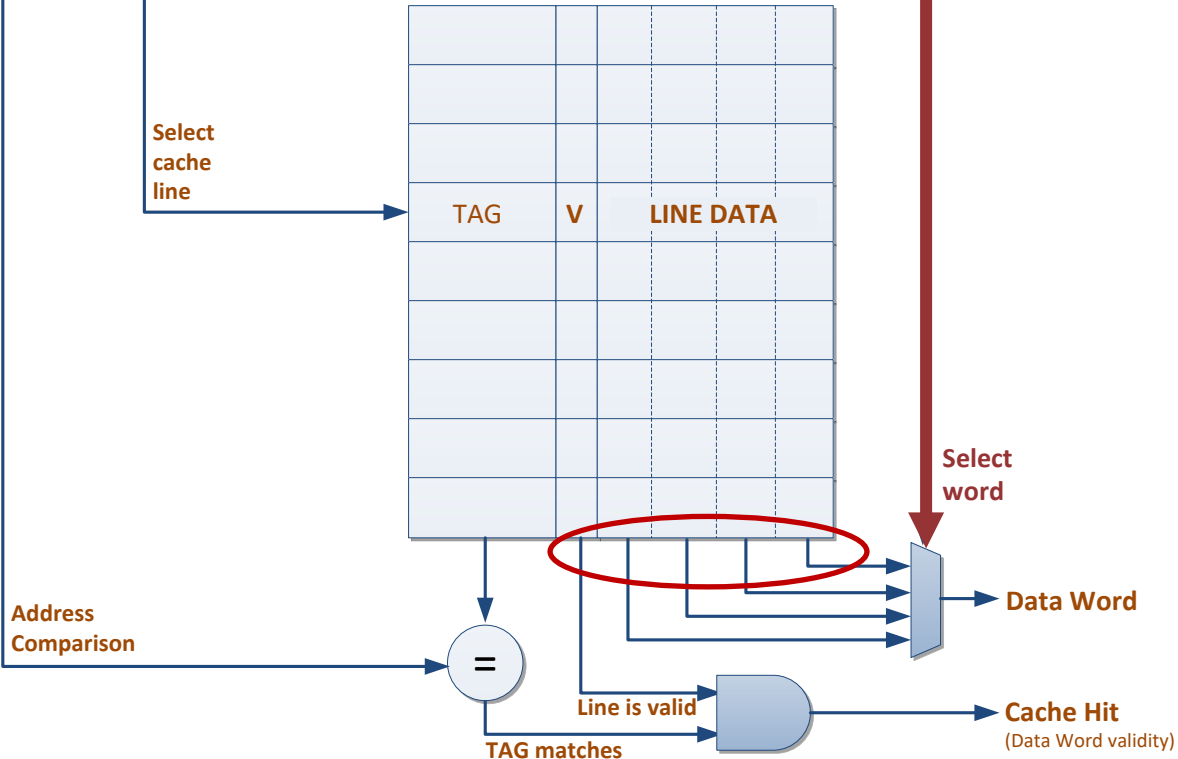
Address



WORD SELECT ou OFFSET

- Seleciona a palavra (dentro do bloco) pretendida, usado apenas na via que deu HIT
- $\#bits(offset) = \log_2(\#palavras)$
- $\#palavras = \text{número de palavras de memória que cabem num único bloco/linha da cache}$

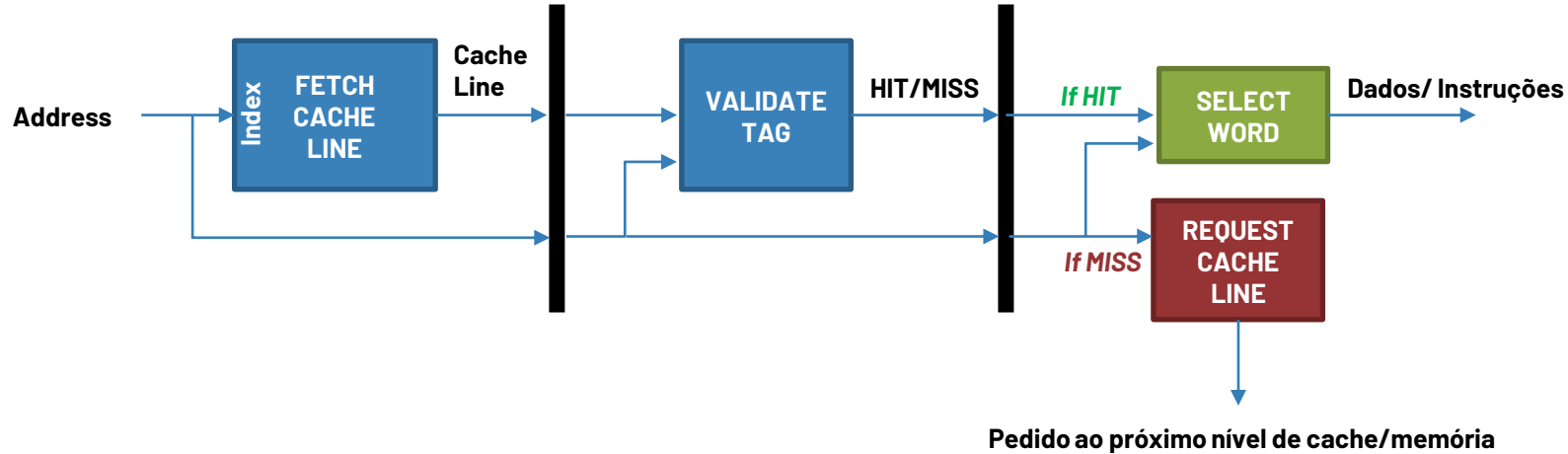
CACHE MEMORY



Acesso em Pipeline

(com 3 estágios)

Possível acesso à cache em Pipeline



Decomposição da palavra de endereço

Considere uma cache L1
para um processador
com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para
32KB
- Linhas com 32B

EXERCÍCIO #1

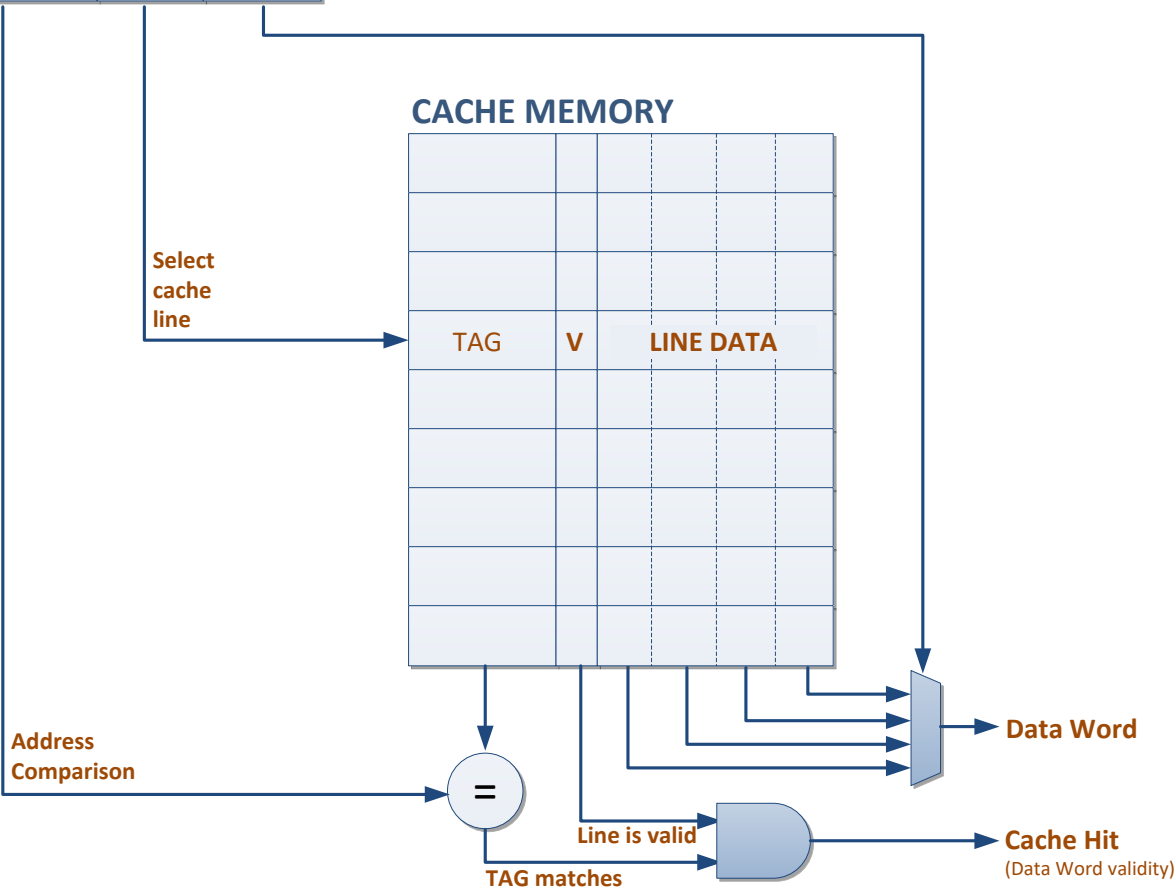
Decomposição da palavra de endereço

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Address

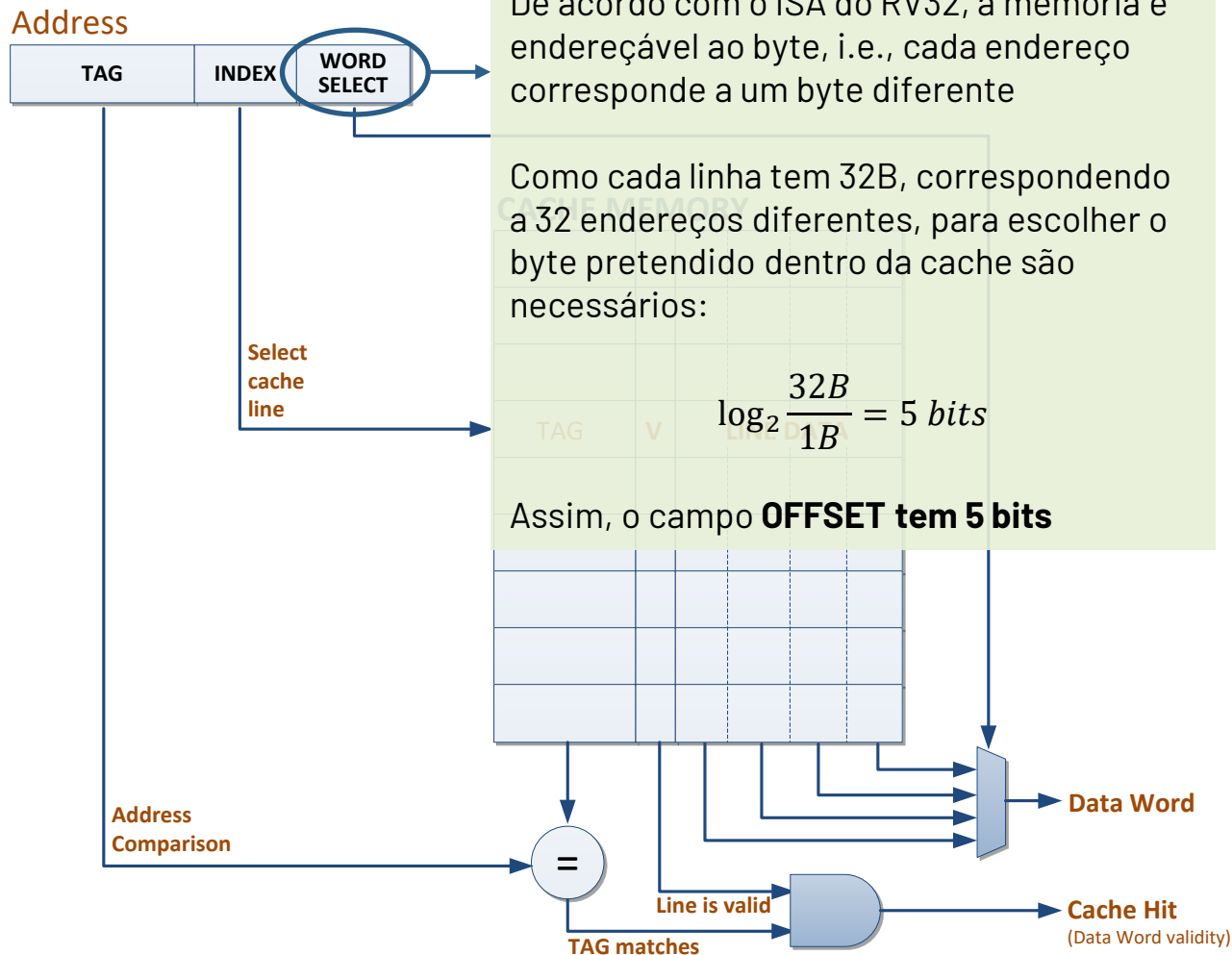


Decomposição da palavra de endereço

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

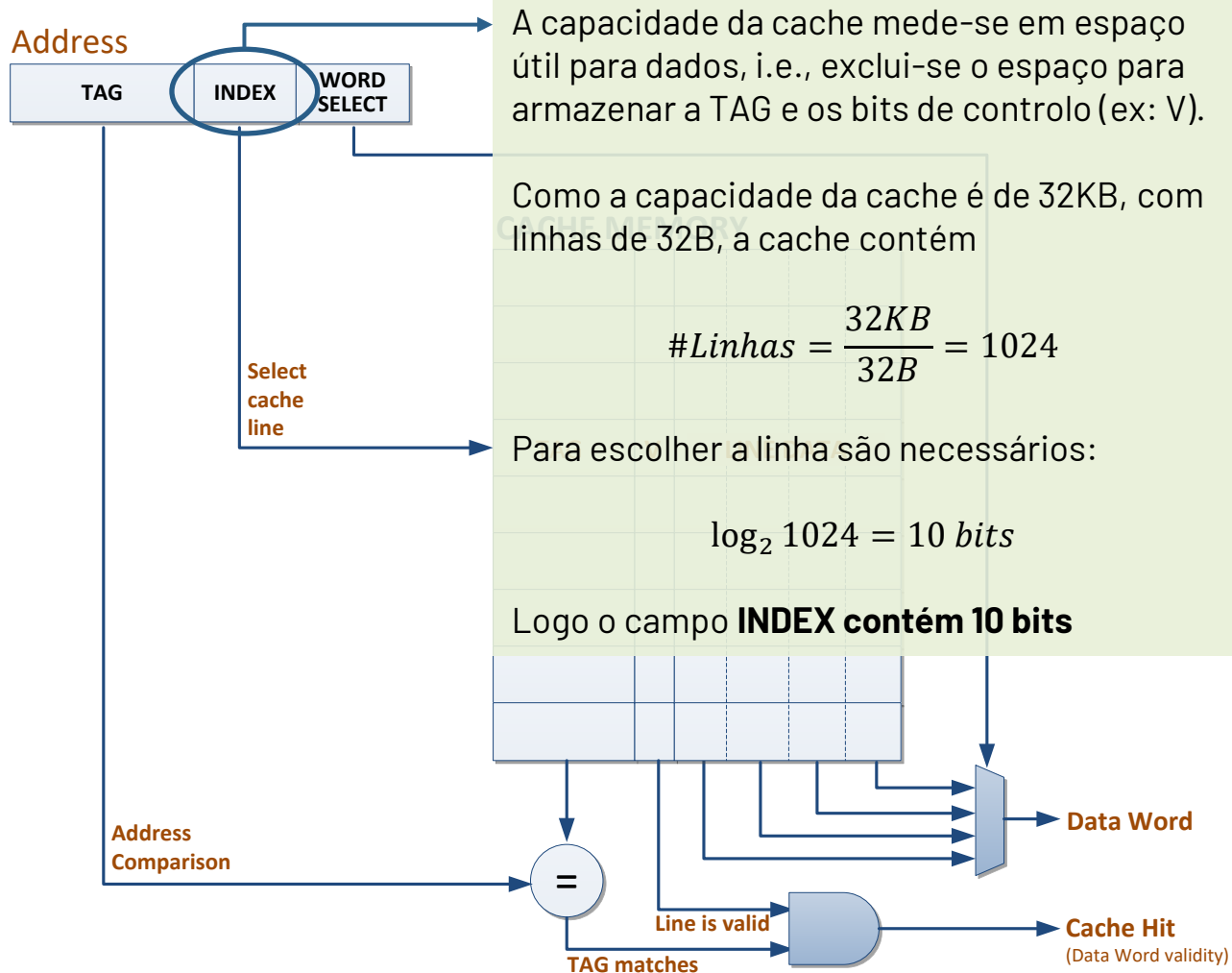


Decomposição da palavra de endereço

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

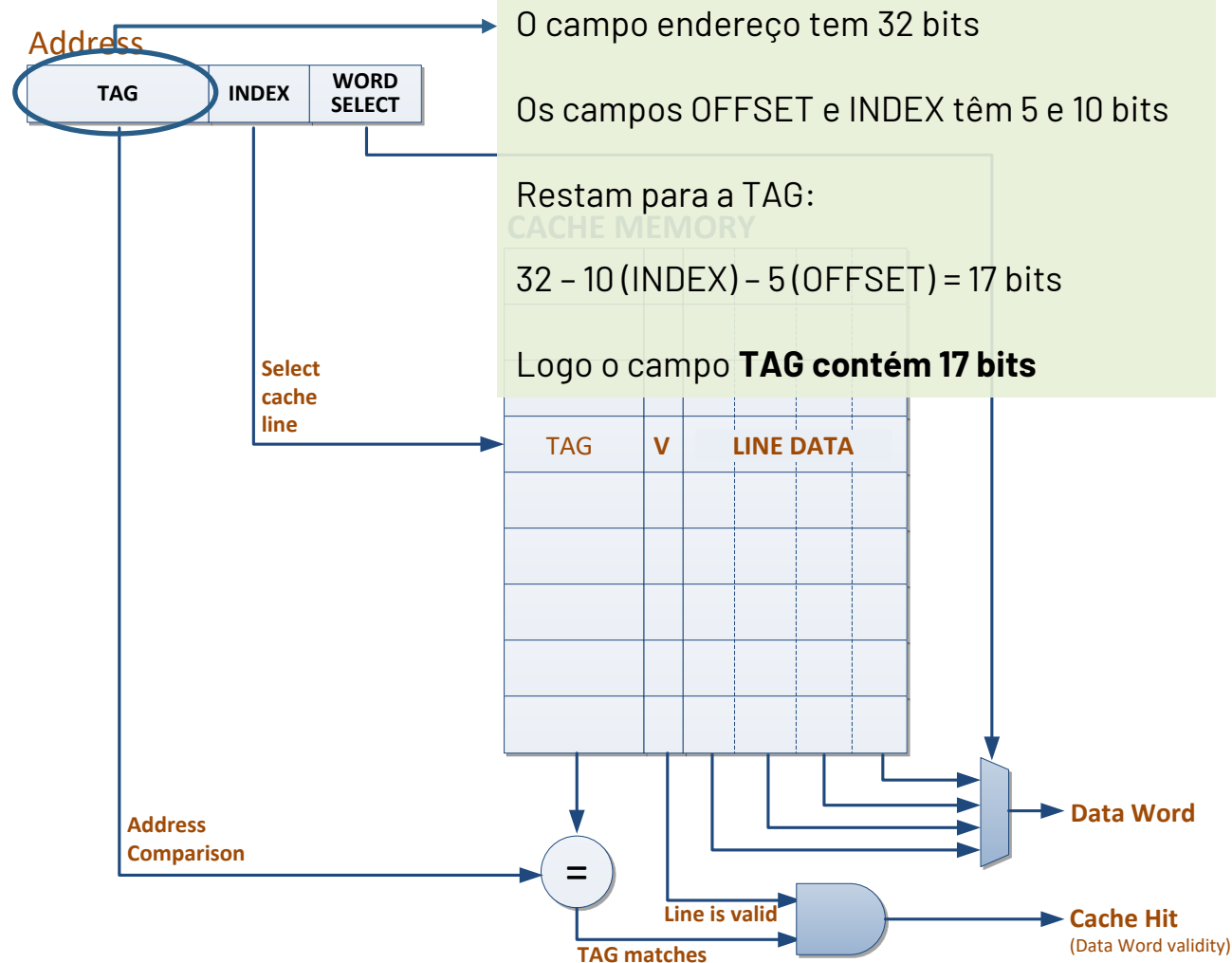


Decomposição da palavra de endereço

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

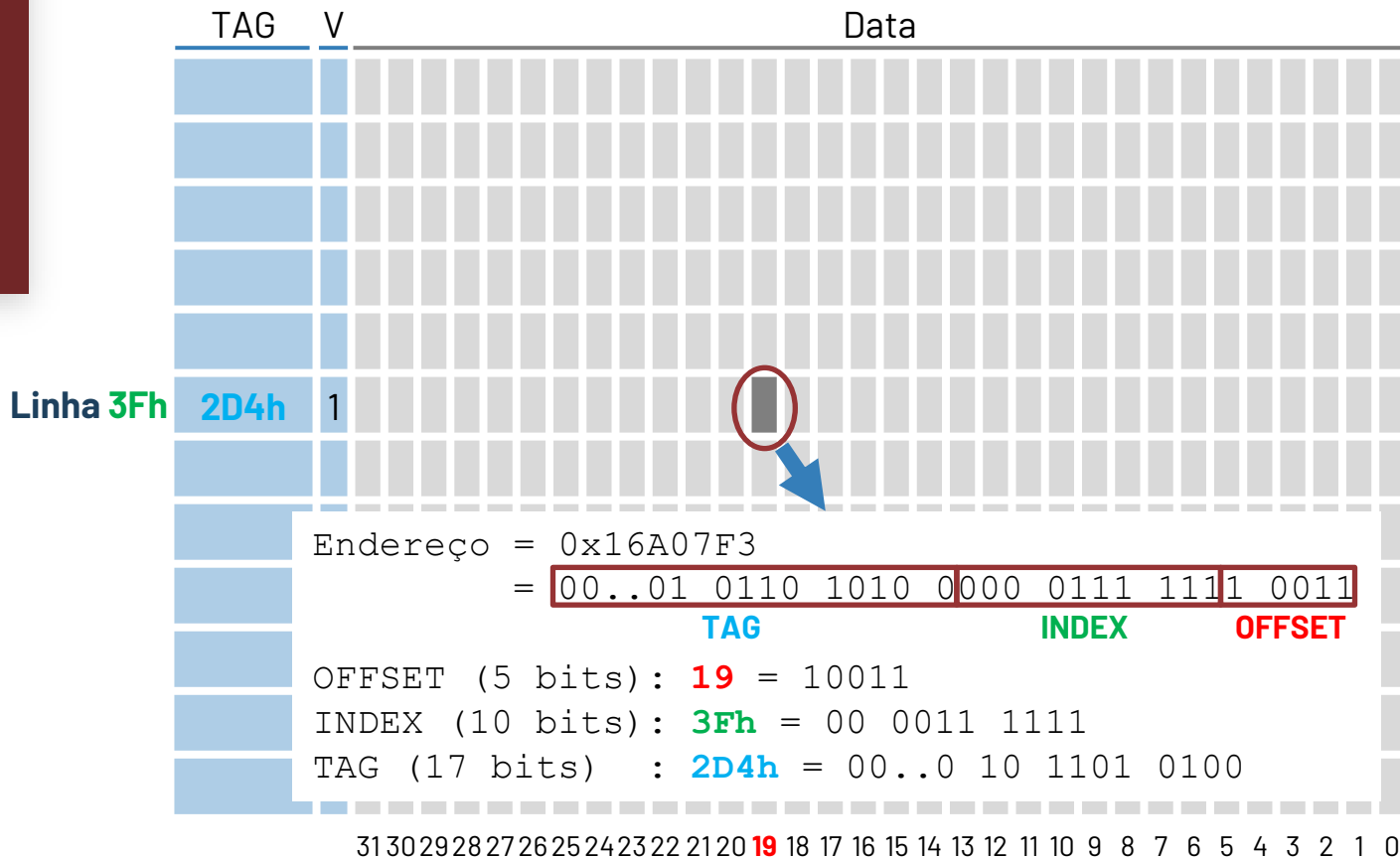


Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B



- Define-se a taxa de sucesso no acesso à cache (HIT RATE) como:

$$\text{HIT RATE} = \frac{\#HITS}{\#ACESSOS}$$

- De forma análoga, define-se o MISS RATE como:

$$\text{MISS RATE} = \frac{\#MISSES}{\#ACESSOS}$$

- Note-se que:

$$\text{MISS RATE} = 1 - \text{HIT RATE}$$

Nota: o hit rate e miss rate são geralmente representados como uma percentagem.

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

EXERCÍCIO #2

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Nota: existem duas caches L1:

- L1-I para instruções
- L1-D para dados

Assim, para o acesso aos dados podemos ignorar o acesso às instruções

```
Vector:  .double  <list of values>

        li        x10, 512
        la        x11, Vector
        addi      x12, x10, -1
        slli      x12, x12, 3
        add       x11, x11, x12
        fcvt.d.w  f0, x0
        ble       x10, x0, fim
loop:    fld       f1, 0(x11)
        fadd.d    f0, f0, f1
        addi      x11, x11, -8
        addi      x10, x10, -1
        bgt       x10, x0, loop
fim:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

X10 = 512

X11 = 0x13F8

X12 = 4088 = 0xFF8

X13 = 0

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop
```

fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Acede ao endereço 0x13F8:

TAG	INDEX	OFFSET
00...00	0001 0011 1111	1 1000

Vamos aceder à linha 9Fh.

Assumindo a cache inicialmente vazia, i.e., o bit de validade (V) é zero em todas as linhas, vamos ter um **MISS!**

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop
```

fim:

Estrutura da cache

Considere uma cache L1
para um processador
com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Estado inicial da cache:

The diagram illustrates the bit fields of a 32-bit memory address. The address is represented as a horizontal bar divided into three sections: TAG (bits 31-24), V (bit 23), and Data (bits 22-0). A callout box titled "Acude ao endereço 0x13F8:" shows the binary representation of this address: TAG=00...00, INDEX=001, and OFFSET=1111. A blue arrow points from the OFFSET field to bit 23 of the address, which is labeled "9F:".

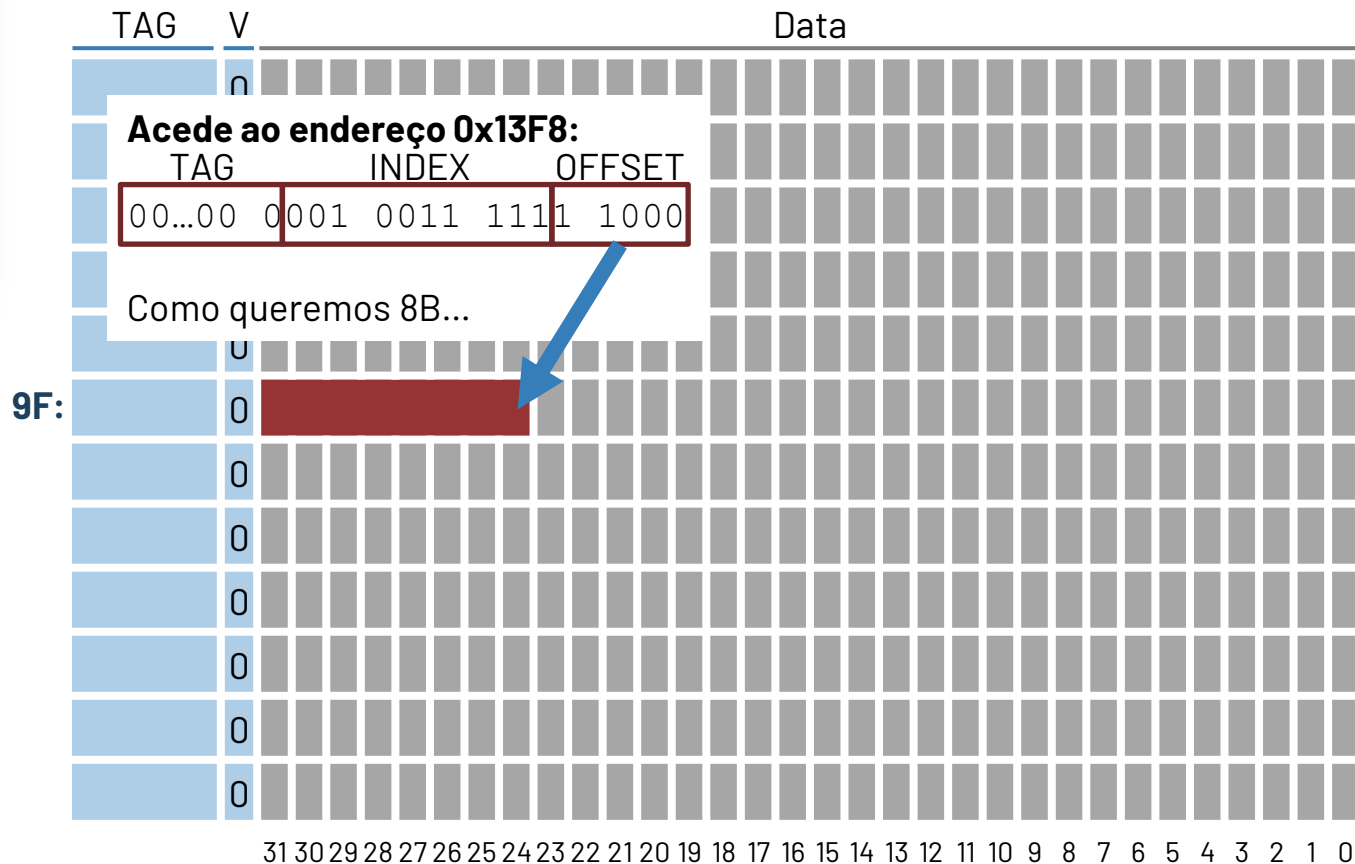
Estrutura da cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Estado inicial da cache:



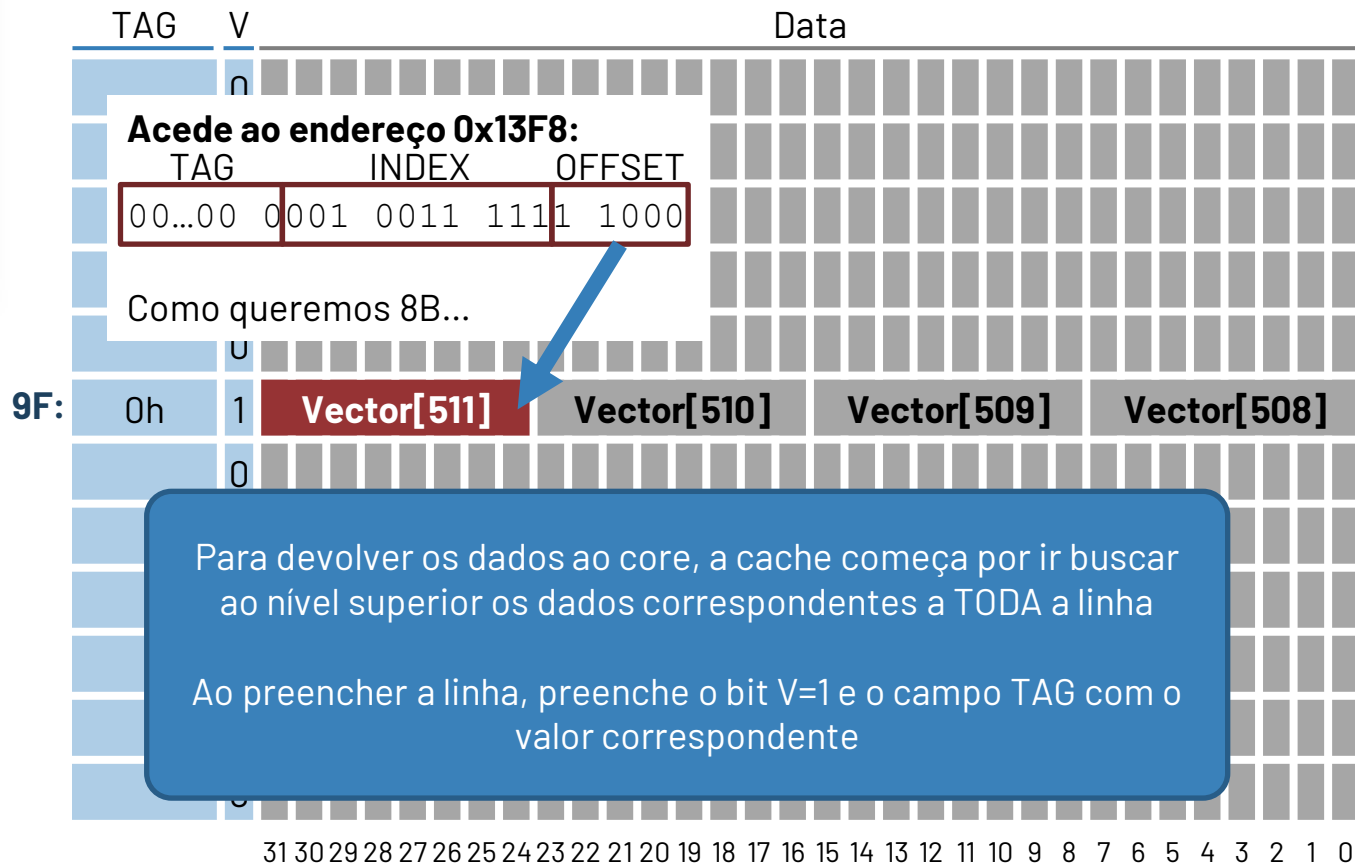
Estrutura da cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Estado da cache:



Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop

fim:
```

X10 = 511
X11 = 0x13F0



fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Acede ao endereço 0x13F0:

TAG	INDEX	OFFSET
00...00	0001 0011 1111	1 0000



Vamos novamente aceder à linha 9Fh.

Como esta entrada foi preenchida com o último acesso, vamos ter um **HIT**.

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop

fim:
```

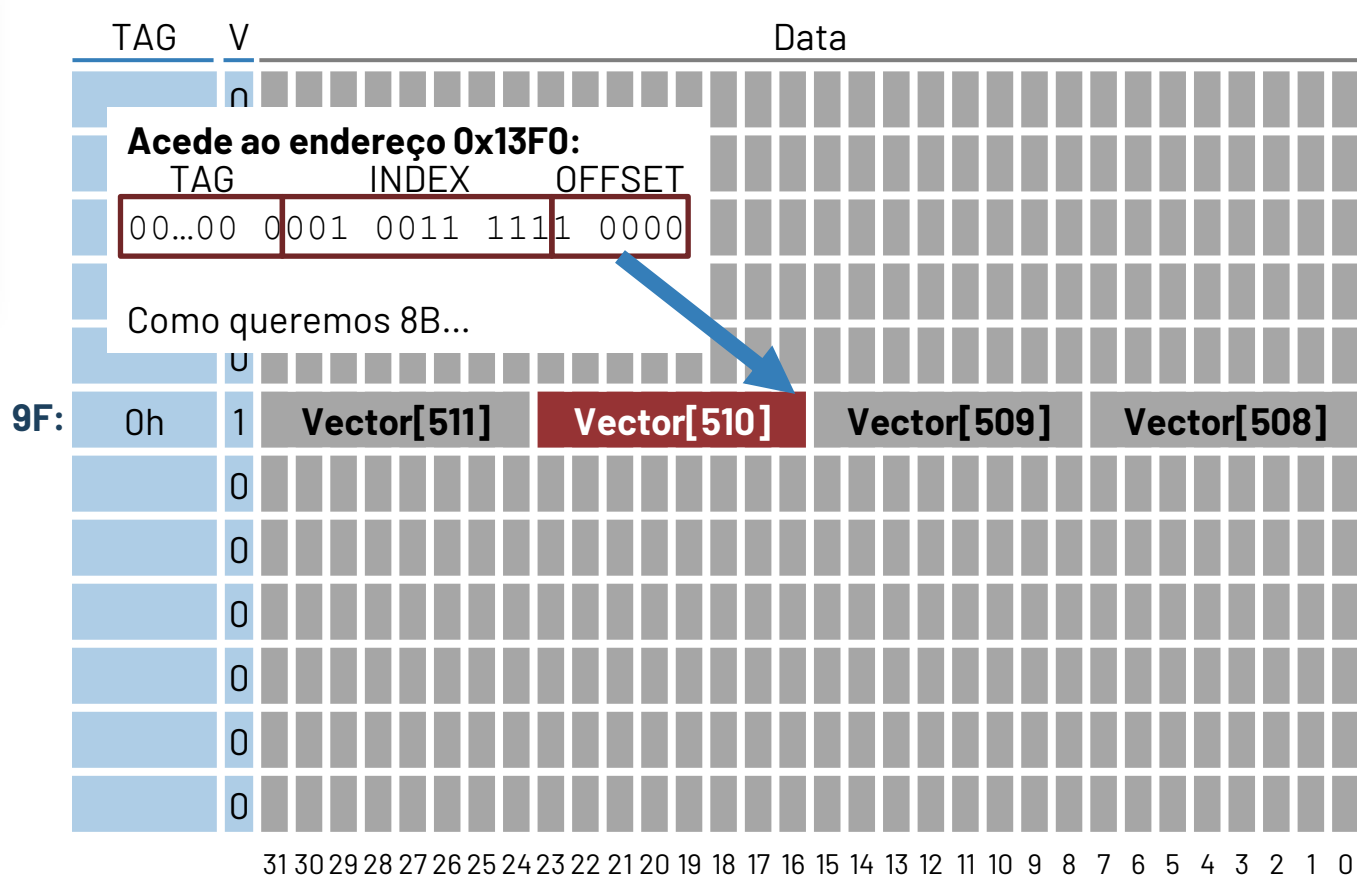
Estrutura da cache

Considere uma cache L1
para um processador
com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Estado da cache:



Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop

fim:
```

X0 = 510
X1 = 0x13E8



fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Acede ao endereço 0x13E8:

TAG	INDEX	OFFSET
00...00	0001 0011 111	0 1000



Vamos aceder novamente à linha 9Fh.

Como esta entrada já foi preenchida vamos ter um **HIT**.

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop

fim:
```


Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop

fim:
```

X10 = 509

X11 = 0x13E0



fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Acede ao endereço 0x13E0:

TAG	INDEX	OFFSET
00...00	0001 0011 1110	0000



Vamos aceder novamente à linha 9Fh.

Como esta entrada já foi preenchida vamos ter um **HIT**.

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop

fim:
```


Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim
loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop
fim:
```

X10 = 508

X11 = 0x13D8



Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Acede ao endereço 0x13D8:

TAG	INDEX	OFFSET
00...00	0001 0011 1101	1 1000



Vamos aceder a uma linha diferente, i.e., linha 9Eh.

Como esta entrada ainda não foi preenchida vamos ter um **MISS**.

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop
```

fim:

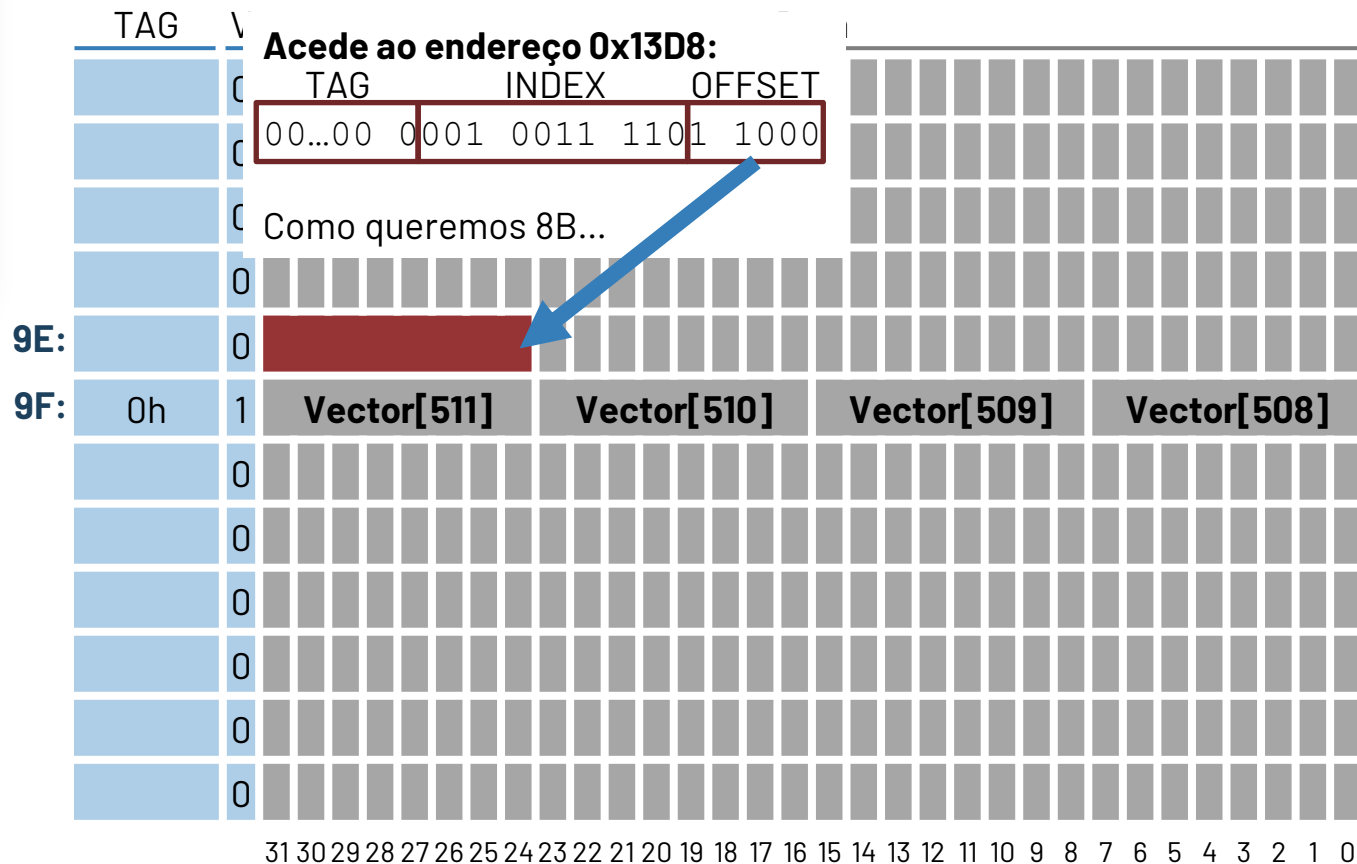
Estrutura da cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Estado da cache antes do acesso ao endereço 13D8h:



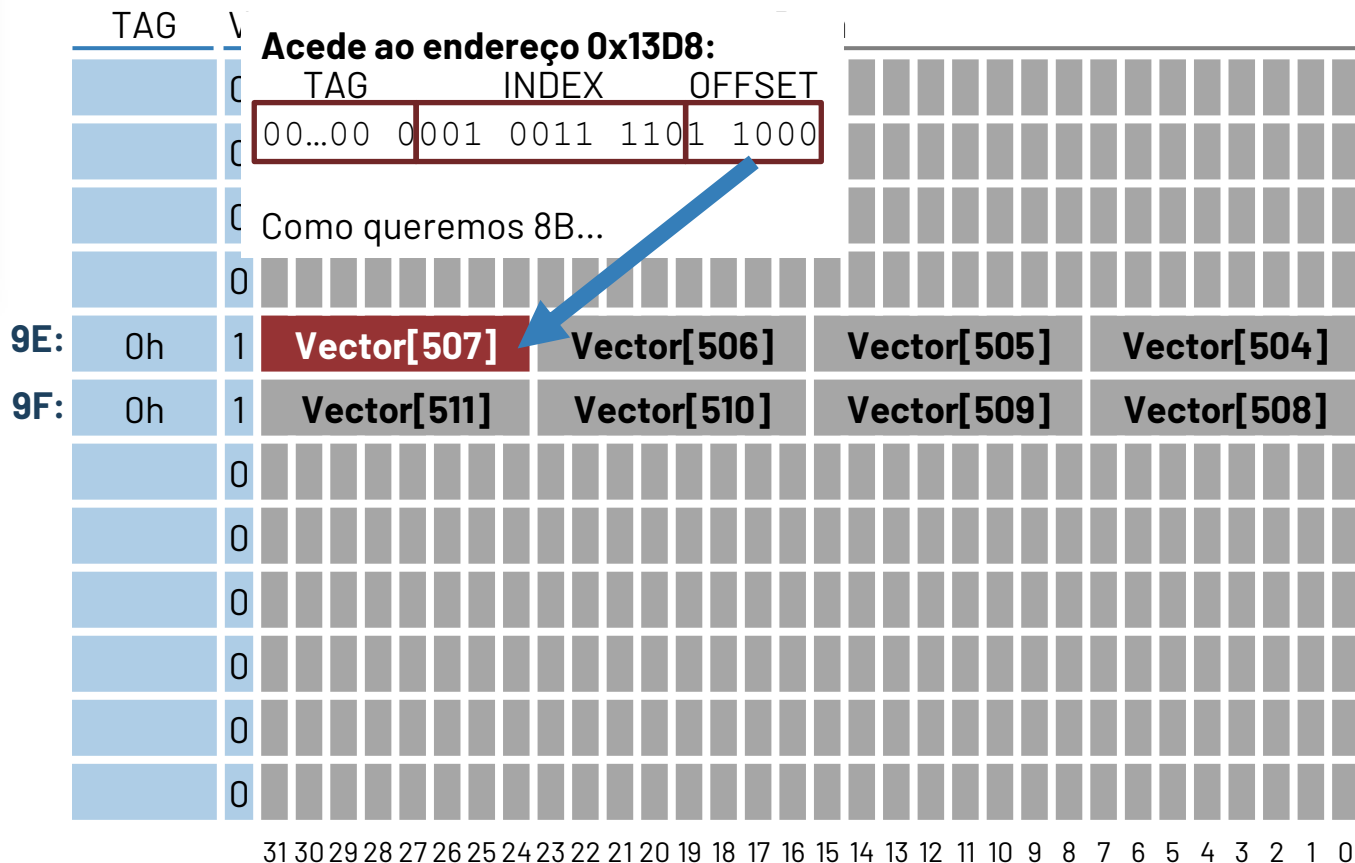
Estrutura da cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Estado da cache após o acesso ao endereço 13D8h:



Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

Para a resolução do exercício, admita que o primeiro elemento do vetor encontra-se no endereço 400h, i.e., Vector=0x400

Acesso aos endereço 0x13D0, 0x13C8, 0x13C0:

TAG	INDEX	OFFSET
00...00	0001 0011 1101	1 1000

TAG	INDEX	OFFSET
00...00	0001 0011 1101	1 1000

TAG	INDEX	OFFSET
00...00	0001 0011 1101	1 1000

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim

loop:   fld      f1, 0(x11)
        fadd.d   f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop
```

fim:

Vamos ter **HIT** (sempre a mesma linha)

Estrutura da cache

Considere uma cache L1
para um processador
com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Estado da cache após o acesso ao endereço 13D8h:

TAG		V	Data																															
		0																																
		0																																
		0																																
		0																																
9E:	0h	1	Vector[507]								Vector[506]								Vector[505]								Vector[504]							
9F:	0h	1	Vector[511]								Vector[510]								Vector[509]								Vector[508]							
		0																																
		0																																
		0																																
		0																																
		0																																
		0																																
		0																																

Acesso aos endereços 0x13D0, 0x13C8, 0x13C0:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.

A taxa de sucesso (HIT RATE):

$$HIT\ RATE = \frac{\#HITS}{\#ACESSOS}$$

- Acesso a um vetor com 512 posições
- No primeiro acesso temos um **MISS**
- Nos 3 seguintes temos um **HIT**
- A padrão repete-se em múltiplos de 4

$$\#ACESSOS = 512$$

$$\#HITS = 3 \times \frac{512}{4} = 384$$

$$HIT\ RATE\ (HR) = \frac{384}{512} = \frac{3}{4} = 75\%$$

$$MISS\ RATE = 1 - HR = 1/4 = 25\%$$

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
slli    x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim
loop:   fld     f1, 0(x11)
        fadd.d  f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop
fim:
```

Tempo (latência) médio de acesso aos dados

- No acesso aos dados, primeiro consulta-se a cache, se houver um MISS consultamos o nível seguinte.
- Para um sistema com apenas um nível de cache obtém-se:

$$T_{Acesso} = T_{L1} + MR_{L1} \times T_{RAM}$$

Onde

MR_{L1} corresponde ao MISS RATE da cache (L1)

Tempo (latência) médio de acesso aos dados

Sistemas hierárquicos com vários níveis de cache

- No acesso aos dados, primeiro consulta-se a cache, se houver um MISS consultamos o nível seguinte.
- Para um sistema com 2 níveis de cache:

$$T_{Acesso} = T_{L1} + MR_{L1} \times (T_{L2} + MR_{L2} \times T_{RAM})$$

Onde

MR_{L1}, MR_{L2} corresponde ao MISS RATE das caches L1 e L2

Tempo (latência) médio de acesso aos dados

Sistemas hierárquicos com vários níveis de cache

- No acesso aos dados, primeiro consulta-se a cache, se houver um MISS consultamos o nível seguinte.
- Para um sistema com 3 níveis de cache:

$$T_{Acesso} = T_{L1} + MR_{L1} \times (T_{L2} + MR_{L2} \times (T_{L3} + MR_{L3} \times T_{RAM}))$$

Onde

$MR_{L1}, MR_{L2}, MR_{L3}$ corresponde ao MISS RATE das caches L1, L2 e L3

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

EXERCÍCIO #3

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine o tempo médio de acesso aos dados

Para a resolução do exercício, admita:

1. Um sistema com apenas 1 nível de cache
2. A latência de acesso à cache L1 é de 2,5 ns
3. A latência de acesso à memória principal é de 100 ns
4. O processador opera a uma frequência de 1GHz

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
sll     x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim
loop:   fld     f1, 0(x11)
        fadd.d  f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop
fim:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine o tempo médio de acesso aos dados

Para a resolução do exercício, admita:

1. Um sistema com apenas 1 nível de cache
2. A latência de acesso à cache L1 é de 2,5 ns
3. A latência de acesso à memória principal é de 100 ns
4. O processador opera a uma frequência de 1GHz

$$T_{CLK} = \frac{1}{1GHz} = 1ns$$

$$Lat_{L1} = \left\lceil \frac{2,5ns}{T_{CLK}} \right\rceil = \left\lceil \frac{2,5ns}{1ns} \right\rceil = 3 \text{ ciclos}$$

$$Lat_{RAM} = \left\lceil \frac{100ns}{T_{CLK}} \right\rceil = \left\lceil \frac{100ns}{1ns} \right\rceil = 100 \text{ ciclos}$$

```
Vector:  .double  <list of values>

li      x10,512
la      x11,Vector
addi    x12,x10,-1
sll     x12,x12,3
add     x11,x11,x12
fcvt.d.w f0,x0
ble     x10,x0,fim

loop:   fld     f1,0(x11)
fadd.d  f0,f0,f1
addi    x11,x11,-8
addi    x10,x10,-1
bgt     x10,x0,loop

fim:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine o tempo médio de acesso aos dados

Para a resolução do exercício, admita:

1. Um sistema com apenas 1 nível de cache
2. A latência de acesso à cache L1 é de 2,5 ns
3. A latência de acesso à memória principal é de 100 ns
4. O processador opera a uma frequência de 1GHz

$$Lat_{Acesso} = Lat_{L1} + MR_{L1} \times Lat_{RAM}$$

$$Lat_{Acesso} = 3 + 0.25 \times 100 = 28 \text{ ciclos}$$

$$T_{Acesso} = Lat_{Acesso} \times T_{CLK} = 28 \times 1ns = 28ns$$

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
sll     x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim
loop:   fld     f1, 0(x11)
        fadd.d  f0, f0, f1
        addi    x11, x11, -8
        addi    x10, x10, -1
        bgt     x10, x0, loop
fim:
```


Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine o tempo médio de acesso aos dados

Para a resolução do exercício, admita:

1. Um sistema com apenas 1 nível de cache
2. A latência de acesso à cache L1 é de 2,5 ns
3. A latência de acesso à memória principal é de 100 ns
4. O processador opera a uma frequência de 1GHz

COM cache:

$$T_{Acesso} = Lat_{Acesso} \times T_{CLK} = 28 \times 1ns = 28ns$$

SEM cache (só com mem. principal):

$$T_{Acesso} = Lat_{Acesso} \times T_{CLK} = 100 \times 1ns = 100ns$$

ou seja: neste caso, a cache proporciona uma aceleração de cerca de $100/28=3,57$ vezes!!!

Vector: .double <list of values>

```
li      x10, 512
la      x11, Vector
addi    x12, x10, -1
sll     x12, x12, 3
add     x11, x11, x12
fcvt.d.w f0, x0
ble     x10, x0, fim
loop:   fld     f1, 0(x11)
fadd.d  f0, f0, f1
addi    x11, x11, -8
addi    x10, x10, -1
bgt     x10, x0, loop
```

fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

EXERCÍCIO #4

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine se existe alguma vantagem em introduzir uma cache L2, com uma latência de 10 ciclos de relógio.

Considere que o miss rate na cache L2 é de 50%.

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine se existe alguma vantagem em introduzir uma cache L2, com uma latência de 10 ciclos de relógio.

Considere que o miss rate na cache L2 é de 50%.

$$\begin{aligned} Lat_{Acesso} &= Lat_{L1} + MR_{L1} \times (Lat_{L2} + MR_{L2} \times Lat_{RAM}) \\ &= 3 + 0.25 \times (10 + 0.5 \times 100) = 18 \text{ ciclos} \end{aligned}$$

No acesso aos dados perdem-se (em media) menos 10 ciclos!

Nota: para um acesso individual, sempre que há um miss, perdem-se mais 10 ciclos. Contudo, em média, a introdução da cache é vantajosa, já que:

- 50% das vezes dá miss, perdem-se mais 10 ciclos ☹
- 50% das vezes dá hit, perdem-se menos 90 ciclos ☺

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

EXERCÍCIO #5

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso às instruções para o troço de código indicado.

```
N:      .word      512
Vector: .double    <list of values>

        lw         x10,N
        la         x11,Vector
        addi       x12,x10,-1
        sll        x12,x12,3
        add        x11,x11,x12
        fcvt.d.w   f0,x0
        ble        x10,x0,fim
loop:    fld        f1,0(x11)
        fadd.d     f0,f0,f1
        addi       x11,x11,-8
        addi       x10,x10,-1
        bgt        x10,x0,loop

fim:
```

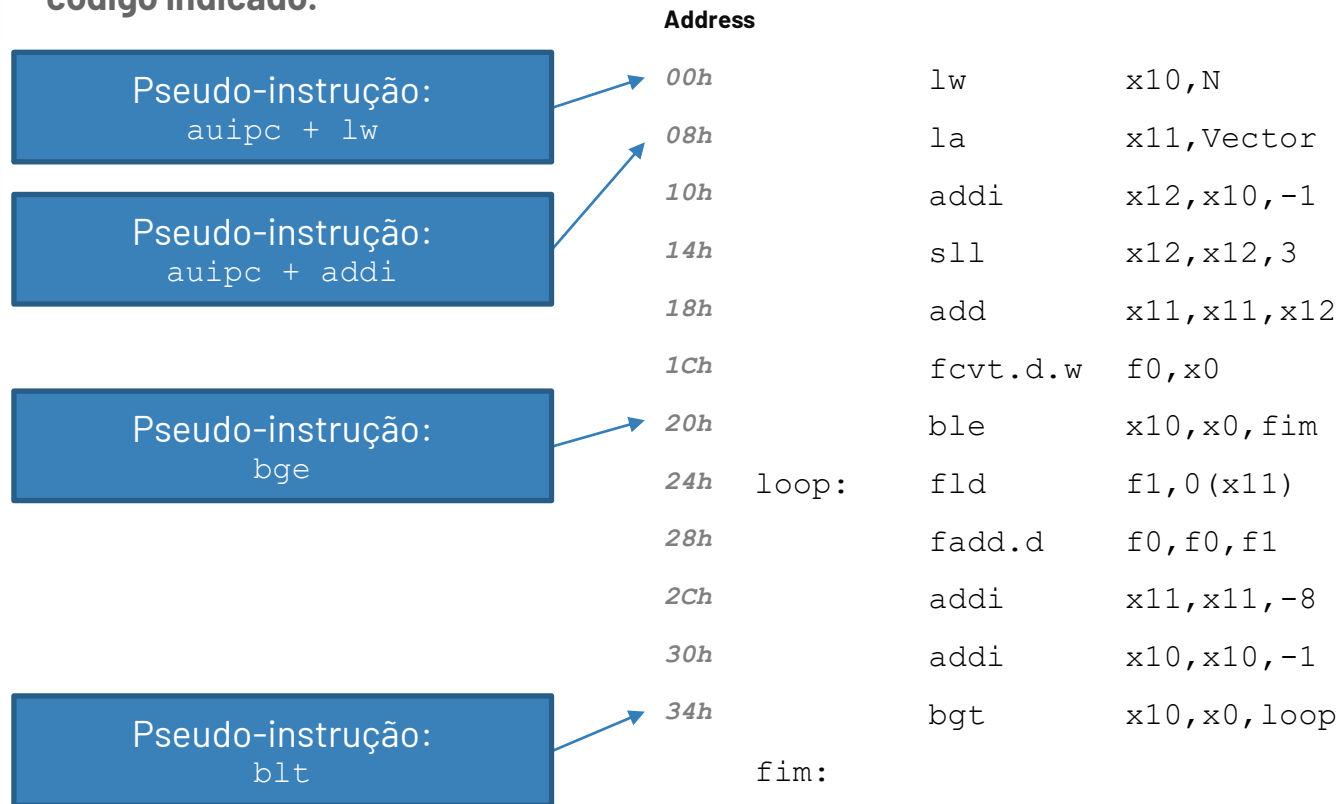
Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso às instruções para o troço de código indicado.



Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso às instruções para o troço de código indicado.

Acede ao endereço 00h:

TAG	INDEX	OFFSET
00...00	0000 0000 0000	0000

Assumindo a cache inicialmente vazia, vamos ter um **MISS!**

Preenchemos a linha 0 da cache...

```
N:      .word      512
Vector: .double    <list of values>
```

Address

```

00h      lw        x10,N
08h
10h      la        x11,Vector
14h      addi      x12,x10,-1
18h      sll       x12,x12,3
1Ch      add       x11,x11,x12
20h      fcvt.d.w  f0,x0
24h      ble      x10,x0,fim
28h      fld       f1,0(x11)
2Ch      fadd.d    f0,f0,f1
30h      addi      x11,x11,-8
34h      addi      x10,x10,-1
        bgt       x10,x0,loop
        fim:
```


Estrutura da cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Estado da cache L1-I após o primeiro acesso:

	TAG	V	Data																															
			la x11,Vector				lw x10,N																											
00h:	0h	1	fcvt(...)	add(...)	sll(...)	addi(...)	addi(...)	auipc(...)	lw(...)	auipc(...)																								
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																																	
	0																													</				

31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso às instruções para o troço de código indicado.

Os próximos 7 acessos vão dar HIT

		N:	.word	512
		Vector:	.double	<list of values>
		Address		
MISS / HIT	00h	lw	x10,N	
HIT / HIT	08h	la	x11,Vector	
HIT	10h	addi	x12,x10,-1	
HIT	14h	sll	x12,x12,3	
HIT	18h	add	x11,x11,x12	
HIT	1Ch	fcvt.d.w	f0,x0	
	20h	ble	x10,x0,fim	
	24h	loop:	fld	f1,0(x11)
	28h		fadd.d	f0,f0,f1
	2Ch		addi	x11,x11,-8
	30h		addi	x10,x10,-1
	34h		bgt	x10,x0,loop
		fim:		

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso às instruções para o troço de código indicado.

Acede ao endereço 20h:

TAG	INDEX	OFFSET
00...00	0000 0000 0010	0000

Corresponde a uma linha diferente

Como a linha ainda não foi preenchida vai dar **MISS**.

Vamos preencher a linha 1 da cache...

```
N:      .word      512
Vector: .double    <list of values>
```

Address

```
00h      lw        x10,N
08h      la        x11,Vector
10h      addi       x12,x10,-1
14h      sll       x12,x12,3
18h      add       x11,x11,x12
1Ch      fcvt.d.w  f0,x0
20h      ble       x10,x0,fim
24h      loop:     fld      f1,0(x11)
28h      fadd.d    f0,f0,f1
2Ch      addi      x11,x11,-8
30h      addi      x10,x10,-1
34h      bgt      x10,x0,loop

fim:
```


Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso às instruções para o troço de código indicado.

Os próximos acessos vão dar **HIT**

Como a cache já está preenchida as próximas iterações do ciclo vão dar sempre **HIT**

Nota: isto só acontece porque nestes exemplos as linhas nunca saem da cache

N: .word 512
Vector: .double <list of values>

		Address		
MISS / HIT	00h	lw	x10,N	
HIT / HIT	08h	la	x11,Vector	
HIT	10h	addi	x12,x10,-1	
HIT	14h	sll	x12,x12,3	
HIT	18h	add	x11,x11,x12	
HIT	1Ch	fcvt.d.w	f0,x0	
MISS	20h	ble	x10,x0,fim	
HIT	24h	loop: fld	f1,0(x11)	
HIT	28h	fadd.d	f0,f0,f1	
HIT	2Ch	addi	x11,x11,-8	
HIT	30h	addi	x10,x10,-1	
HIT	34h	bgt	x10,x0,loop	
		fim:		

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 32KB
- Linhas com 32B

Determine a taxa de sucesso no acesso às instruções para o troço de código indicado.

Preâmbulo:

- 2 MISS, 7 HITs, 9 Acessos

Ciclo:

- 5 HITs, 5 Acessos

Total:

- 2 MISS
- 5x512 HITs
- 2+5x512 Acessos

$$HIT\ RATE\ (HR) = \frac{2567}{2569} = 99.92\%$$

$$MISS\ RATE\ (MR) = \frac{2}{2569} = 0.08\% = 1 - HR$$

N: .word 512
Vector: .double <list of values>

Address

MISS / HIT	00h	lw	x10,N
HIT / HIT	08h	la	x11,Vector
HIT	10h	addi	x12,x10,-1
HIT	14h	sll	x12,x12,3
HIT	18h	add	x11,x11,x12
HIT	1Ch	fcvt.d.w	f0,x0
MISS	20h	ble	x10,x0,fim
HIT	24h	loop: fld	f1,0(x11)
HIT	28h	fadd.d	f0,f0,f1
HIT	2Ch	addi	x11,x11,-8
HIT	30h	addi	x10,x10,-1
HIT	34h	bgt	x10,x0,loop
		fim:	

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

EXERCÍCIO #6

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

Admita os seguintes endereços para a localização dos vetores VA e VB:

VA = &(VA[0]) = 1000h
VB = &(VB[0]) = 2000h

E ainda as seguintes latências de acessos aos dados:

Latência da L1 = 4 ciclos
Latência da DRAM = 50 ciclos

1. **Determine a taxa de sucesso no acesso aos dados para o troço de código indicado.**
2. **Determine o tempo médio de acesso aos dados**

VA: .word <list of values>
VB: .word <list of values>

```
li      x10, 1024 #N
la      x11, VA
la      x12, VB
addi    x13, x10, -1
sll     x13, x13, 2
add     x14, x11, x13
add     x15, x12, x13
mv      x16, x0
ble     x10, x0, fim
loop:   lw      x17, 0(x14)
        lw      x18, 0(x15)
        mul     x17, x17, x18
        add     x16, x16, x17
        addi    x14, x14, -4
        addi    x15, x15, -4
        addi    x10, x10, -1
        bgt     x10, x0, loop
```

fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

Bloco de dados de 16B, i.e., 16 endereços de 1 byte. Para seleccionar o byte, precisamos de:

$$\log_2 \frac{16B}{1B} = 4 \text{ bits}$$

OFFSET → 4 bits

```
VA:      .word  <list of values>
VB:      .word  <list of values>
```

```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim
loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

4KB de dados, distribuídos em linhas de 16B, o que corresponde a $4\text{KB}/16\text{B} = 2^{12}/2^4 = 2^8 = 256$ linhas

Para endereçar 256 linhas são necessários

$$\log_2(256) = 8 \text{ bits}$$

INDEX → 8 bits

```
VA:      .word  <list of values>
VB:      .word  <list of values>
```

```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim
loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

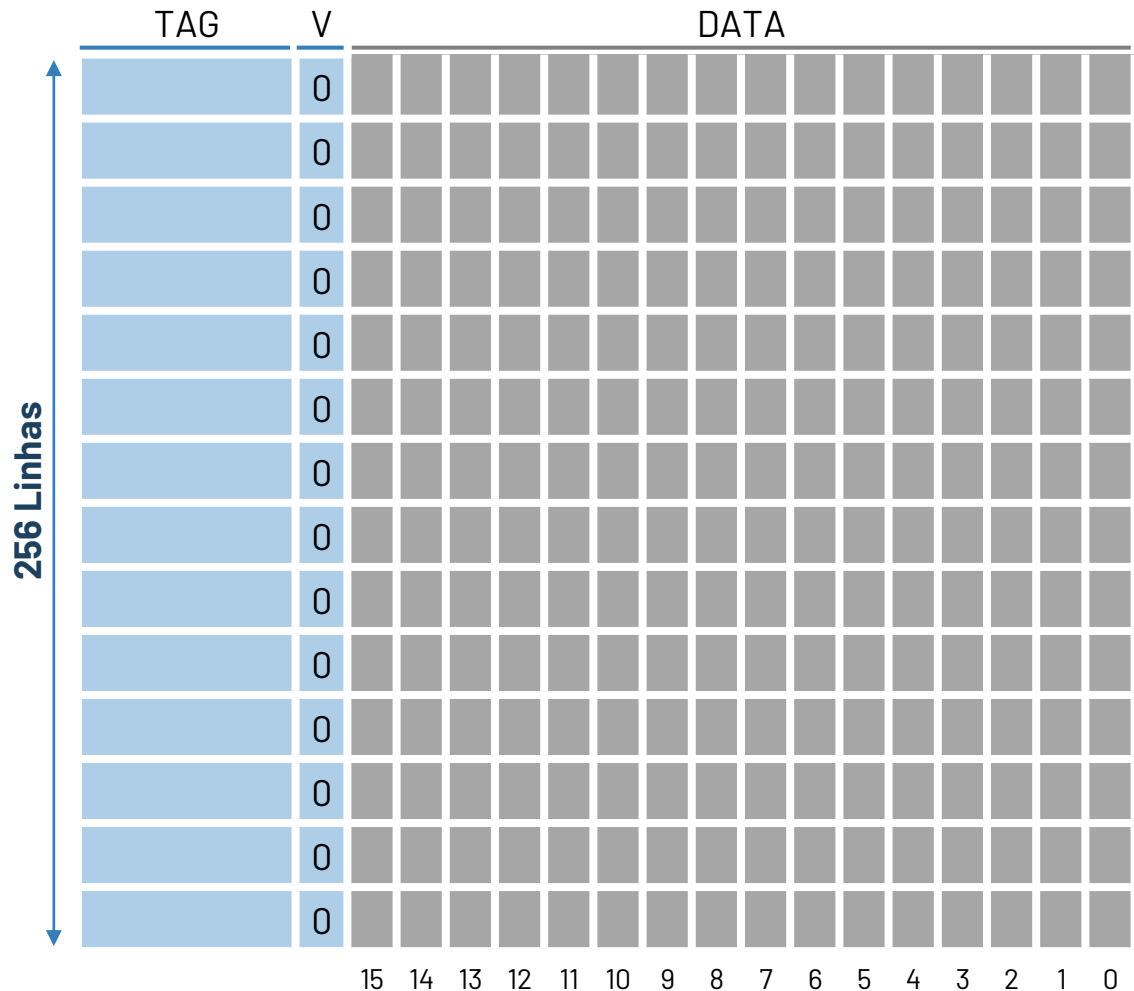
fim:

Estrutura da cache

ESTADO INICIAL DA CACHE:

ADDRESS (32 bits)

TAG	INDEX	OFFSET
20 bits	8 bits	4 bits



Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

Admita os seguintes endereços para a localização dos vetores VA e VB:

VA = &(VA[0]) = 1000h

VB = &(VB[0]) = 2000h

X10 = 1024
X11 = 0x1000
X12 = 0x2000
X13 = 0xFFC
X14 = 0x1FFC
X15 = 0x2FFC
X16 = 0

VA: .word <list of values>
VB: .word <list of values>

```
li      x10, 1024 #N
la      x11, VA
la      x12, VB
addi    x13, x10, -1
sll     x13, x13, 2
add     x14, x11, x13
add     x15, x12, x13
mv      x16, x0
ble     x10, x0, fim
loop:   lw      x17, 0(x14)
        lw      x18, 0(x15)
        mul     x17, x17, x18
        add     x16, x16, x17
        addi    x14, x14, -4
        addi    x15, x15, -4
        addi    x10, x10, -1
        bgt     x10, x0, loop
```

fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

ADDRESS (32 bits)

0 0001h	FFh	Ch
20 bits	8 bits	4 bits

VA: .word <list of values>
VB: .word <list of values>

Como a cache está inicialmente vazia, i.e., todas as entradas têm V=0, o acesso à entrada FFh corresponde necessariamente a um **MISS**!

Vamos preencher a linha...

Endereço 1FFCh

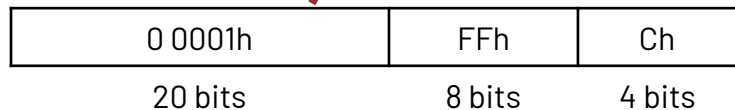
```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim
lw      x17,0(x14)
lw      x18,0(x15)
mul     x17,x17,x18
add     x16,x16,x17
addi    x14,x14,-4
addi    x15,x15,-4
addi    x10,x10,-1
bgt     x10,x0,loop
```

fim:

Estrutura da cache

Preenchimento da cache após o primeiro acesso:

ADDRESS (32 bits)



FFh:

TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
1h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
		ADDRESS: 1FFCh				ADDRESS: 1FF8h				ADDRESS: 1FF4h				ADDRESS: 1FF0h			

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

ADDRESS (32 bits)

0 0002h	FFh	Ch
20 bits	8 bits	4 bits

VA: .word <list of values>
VB: .word <list of values>

Vamos verificar o estado da cache...

```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim

loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

Endereço 2FFCh

fim:

Estrutura da cache

Endereço pedido: 2FFCh

A linha selecionada (FFh), é válida, mas a TAG não é igual. Assim, temos um **MISS**!

ADDRESS (32 bits)

0 0002h	FFh	Ch
20 bits	8 bits	4 bits

A TAG é diferente!

FFh:

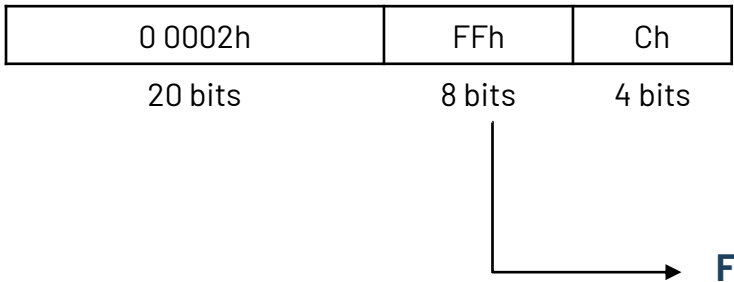
TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
1h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
		ADDRESS: 1FFCh				ADDRESS: 1FF8h				ADDRESS: 1FF4h				ADDRESS: 1FF0h			

Estrutura da cache

Endereço pedido: 2FFCh

Vamos colocar os dados em cache... como corresponde
à mesma linha vamos substituir esta linha...

ADDRESS (32 bits)



TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
2h	1	VB[1023]				VB[1022]				VB[1021]				VB[1020]			
		ADDRESS: 2FFCh				ADDRESS: 2FF8h				ADDRESS: 2FF4h				ADDRESS: 2FF0h			

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

ADDRESS (32 bits)

0 0001h	FFh	8h
20 bits	8 bits	4 bits

VA: .word <list of values>
VB: .word <list of values>

Vamos verificar o estado da cache...

Endereço 1FF8h



```

li      x10, 1024 #N
la      x11, VA
la      x12, VB
addi    x13, x10, -1
sll     x13, x13, 2
add     x14, x11, x13
add     x15, x12, x13
mv      x16, x0
ble     x10, x0, fim
lw      x17, 0(x14)
lw      x18, 0(x15)
mul     x17, x17, x18
add     x16, x16, x17
addi    x14, x14, -4
addi    x15, x15, -4
addi    x10, x10, -1
bgt     x10, x0, loop

```

fim:

Estrutura da cache

Endereço pedido: 1FFCh

A linha selecionada (FFh), é válida, mas a TAG não é igual. Assim, temos um **MISS!**

ADDRESS (32 bits)

0 0001h	FFh	8h
20 bits	8 bits	4 bits

A TAG é diferente!

FFh:

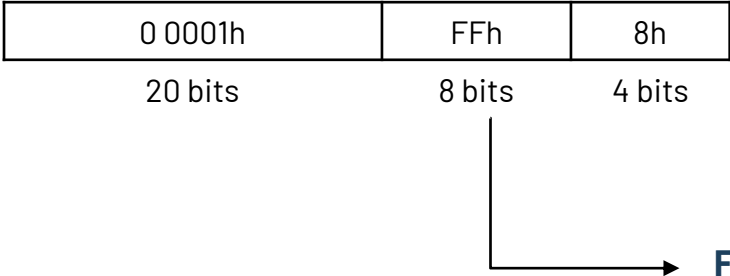
TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
2h	1	VB[1023]				VB[1022]				VB[1021]				VB[1020]			
		ADDRESS: 2FFCh				ADDRESS: 2FF8h				ADDRESS: 2FF4h				ADDRESS: 2FF0h			

Estrutura da cache

Endereço pedido: 1FF8h

Vamos colocar os dados em cache... como corresponde
à mesma linha vamos substituir esta linha...

ADDRESS (32 bits)



TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
1h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
		ADDRESS: 1FFCh				ADDRESS: 1FF8h				ADDRESS: 1FF4h				ADDRESS: 1FF0h			

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

ADDRESS (32 bits)

0 0002h	FFh	8h
20 bits	8 bits	4 bits

VA: .word <list of values>
VB: .word <list of values>

Vamos verificar o estado da cache...

```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim

loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

Endereço 2FF8h

fim:

Estrutura da cache

Endereço pedido: 2FF8h

A linha selecionada (FFh), é válida, mas a TAG não é igual. Assim, temos um **MISS!**

ADDRESS (32 bits)

0 0002h	FFh	8h
20 bits	8 bits	4 bits

A TAG é diferente!

FFh:

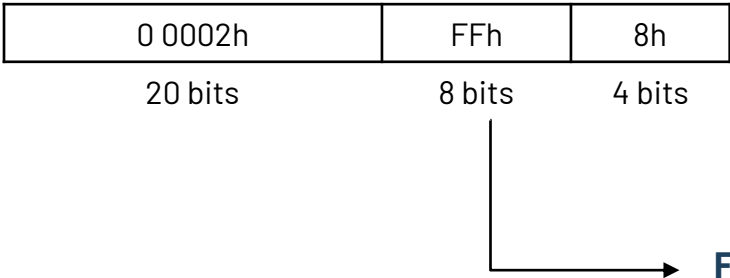
TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
1h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
		ADDRESS: 1FFCh				ADDRESS: 1FF8h				ADDRESS: 1FF4h				ADDRESS: 1FF0h			

Estrutura da cache

Endereço pedido: 2FFCh

Vamos colocar os dados em cache... como corresponde
à mesma linha vamos substituir esta linha...

ADDRESS (32 bits)



TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
	0																
2h	1	VB[1023]				VB[1022]				VB[1021]				VB[1020]			
		ADDRESS: 2FFCh				ADDRESS: 2FF8h				ADDRESS: 2FF4h				ADDRESS: 2FF0h			

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Mapeamento direto
- Capacidade para 4KB
- Linhas com 16B

Embora seja um exemplo feito “à medida”, ilustra um caso típico de padrões de acesso à memória com conflitos.

Em consequência todos os acessos dão **MISS** na cache!

$MISS\ RATE = 100\%$, $HIT\ RATE = 0\%$

Como resolver o problema dos conflitos?

Vamos experimentar uma cache associativa...

```
VA:      .word  <list of values>
VB:      .word  <list of values>
```

```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim
loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- **2 vias (i.e., 2 tabelas)**
- Capacidade para 4KB
- Linhas com 16B

```
VA:      .word  <list of values>
VB:      .word  <list of values>

li       x10,1024 #N
la       x11,VA
la       x12,VB
addi     x13,x10,-1
sll      x13,x13,2
add      x14,x11,x13
add      x15,x12,x13
mv       x16,x0
ble      x10,x0,fim
loop:    lw      x17,0(x14)
         lw      x18,0(x15)
         mul     x17,x17,x18
         add     x16,x16,x17
         addi    x14,x14,-4
         addi    x15,x15,-4
         addi    x10,x10,-1
         bgt     x10,x0,loop

fim:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- 2 vias (i.e., 2 tabelas)
- Capacidade para 4KB
- Linhas com 16B

Bloco de dados de 16B, i.e., 16 endereços de 1 byte. Para seleccionar o byte, precisamos de:

$$\log_2 \frac{16B}{1B} = 4 \text{ bytes}$$

OFFSET → 4 bits

```
VA:      .word  <list of values>
VB:      .word  <list of values>
```

```
li      x10, 1024 #N
la      x11, VA
la      x12, VB
addi    x13, x10, -1
sll     x13, x13, 2
add     x14, x11, x13
add     x15, x12, x13
mv      x16, x0
ble     x10, x0, fim
loop:   lw      x17, 0(x14)
        lw      x18, 0(x15)
        mul     x17, x17, x18
        add     x16, x16, x17
        addi    x14, x14, -4
        addi    x15, x15, -4
        addi    x10, x10, -1
        bgt     x10, x0, loop
```

fim:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- 2 vias (i.e., 2 tabelas)
- Capacidade para 4KB
- Linhas com 16B

4KB de dados, distribuídos em linhas de 16B, o que corresponde a $4\text{KB}/16\text{B}=256$ linhas

Como existem duas vias (tabelas), as linhas distribuem-se de igual forma por ambas.

Assim, cada via tem 128 linhas. Para as endereçar são necessários

$$\log_2(128) = 7 \text{ bits}$$

INDEX → 7 bits

```
VA:      .word  <list of values>
VB:      .word  <list of values>
```

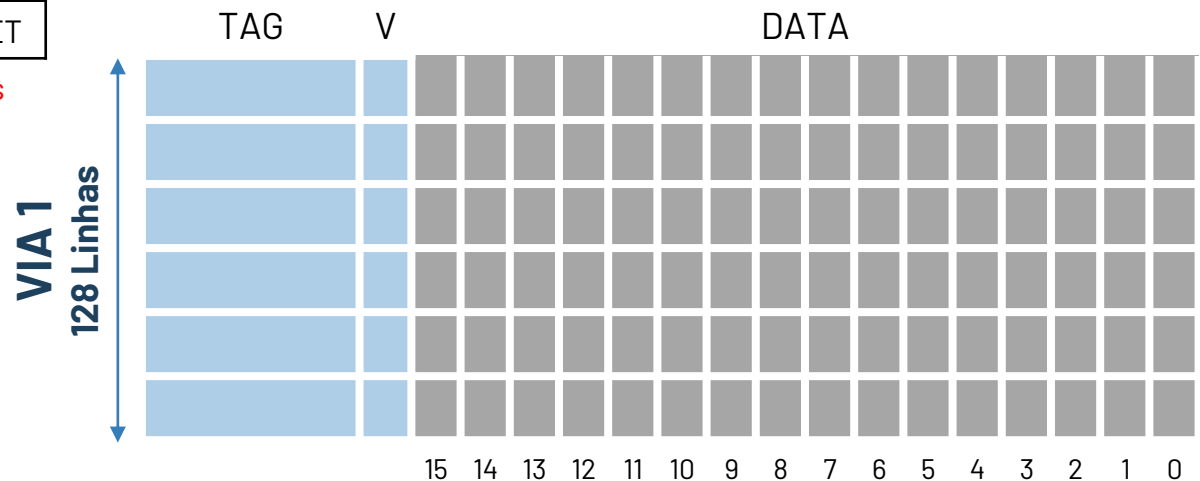
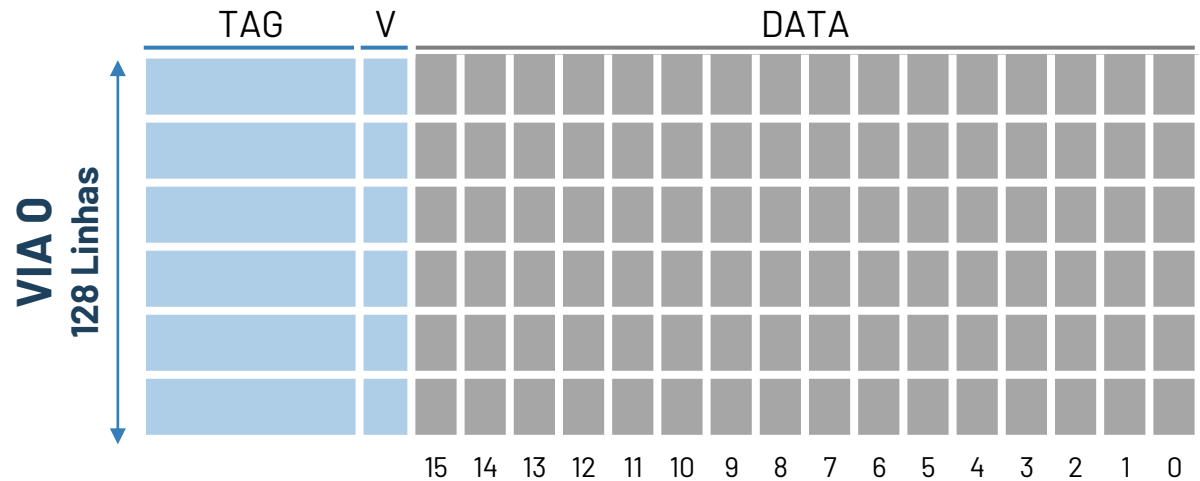
```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim
loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

fim:

Estrutura da cache

ADDRESS (32 bits)

TAG	INDEX	OFFSET
21 bits	7 bits	4 bits



Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- 2 vias (i.e., 2 tabelas)
- Capacidade para 4KB
- Linhas com 16B

ADDRESS (32 bits)

3h	7Fh	Ch
21 bits	7 bits	4 bits

VA: .word <list of values>
VB: .word <list of values>

Como a cache está inicialmente vazia, i.e., todas as entradas têm V=0, o acesso à entrada 7Fh corresponde necessariamente a um **MISS**!

Vamos preencher a linha...

Endereço 1FFCh



```
li      x10, 1024 #N
la      x11, VA
la      x12, VB
addi    x13, x10, -1
sll     x13, x13, 2
add     x14, x11, x13
add     x15, x12, x13
mv      x16, x0
ble     x10, x0, fim
lw      x17, 0(x14)
lw      x18, 0(x15)
mul     x17, x17, x18
add     x16, x16, x17
addi    x14, x14, -4
addi    x15, x15, -4
addi    x10, x10, -1
bgt     x10, x0, loop
```

fim:

Estrutura da cache

ADDRESS (32 bits)

3h	7F	Ch
21 bits	7 bits	4 bits

Endereço pedido: 1FFCh

MISS....

Ambas as vias estão vazias... escolhemos uma delas

VIA 0

TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
7Fh: 3h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
		15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

VIA 1

TAG	V	DATA															
	0																
	0																
	0																
	0																
	0																
	0																
		15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- 2 vias (i.e., 2 tabelas)
- Capacidade para 4KB
- Linhas com 16B

ADDRESS (32 bits)

5h	7Fh	Ch
21 bits	7 bits	4 bits

VA: .word <list of values>
VB: .word <list of values>

Vamos verificar o estado da cache...

```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim

loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

Endereço 2FFCh

fim:

Estrutura da cache

ADDRESS (32 bits)

5h	7F	Ch
21 bits	7 bits	4 bits

Endereço pedido: 2FFCh

MISS....

A via 0 está preenchida, vamos preencher a via 1

VIA 0

	TAG	V	DATA															
		0																
		0																
		0																
		0																
		0																
7Fh:	3h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
			15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

VIA 1

	TAG	V	DATA															
		0																
		0																
		0																
		0																
		0																
7Fh:	5h	1	VB[1023]				VB[1022]				VB[1021]				VB[1020]			
			15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- 2 vias (i.e., 2 tabelas)
- Capacidade para 4KB
- Linhas com 16B

ADDRESS (32 bits)

0 0003h	7Fh	8h
21 bits	7 bits	4 bits

VA: .word <list of values>
VB: .word <list of values>

Vamos verificar o estado da cache...

Endereço 1FF8h



```
li      x10, 1024 #N
la      x11, VA
la      x12, VB
addi    x13, x10, -1
sll     x13, x13, 2
add     x14, x11, x13
add     x15, x12, x13
mv      x16, x0
ble     x10, x0, fim
lw      x17, 0(x14)
lw      x18, 0(x15)
mul     x17, x17, x18
add     x16, x16, x17
addi    x14, x14, -4
addi    x15, x15, -4
addi    x10, x10, -1
bgt     x10, x0, loop
```

fim:

Estrutura da cache

ADDRESS (32 bits)

3h	7F	8h
21 bits	7 bits	4 bits

Endereço pedido: 1FF8h

HIT na via 0

VIA 0

	TAG	V	DATA															
		0																
		0																
		0																
		0																
		0																
7Fh:	3h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
			15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

VIA 1

	TAG	V	DATA															
		0																
		0																
		0																
		0																
		0																
7Fh:	5h	1	VB[1023]				VB[1022]				VB[1021]				VB[1020]			
			15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- 2 vias (i.e., 2 tabelas)
- Capacidade para 4KB
- Linhas com 16B

ADDRESS (32 bits)

5h	7Fh	8h
21 bits	7 bits	4 bits

VA: .word <list of values>
VB: .word <list of values>

Vamos verificar o estado da cache...

```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim

loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

Endereço 2FF8h

fim:

Estrutura da cache

ADDRESS (32 bits)

5h	7F	8h
21 bits	7 bits	4 bits

Endereço pedido: 2FF8h

HIT na via 1

VIA 0

	TAG	V	DATA															
		0																
		0																
		0																
		0																
		0																
7Fh:	3h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
			15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

VIA 1

	TAG	V	DATA															
		0																
		0																
		0																
		0																
		0																
7Fh:	5h	1	VB[1023]				VB[1022]				VB[1021]				VB[1020]			
			15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- 2 vias (i.e., 2 tabelas)
- Capacidade para 4KB
- Linhas com 16B

Próximos acessos:

ADDRESS (32 bits): 1FF4h

3h	7Fh	4h
----	-----	----

HIT na via 0

ADDRESS (32 bits): 2FF4h

5h	7Fh	4h
----	-----	----

HIT na via 1

ADDRESS (32 bits): 1FF0h

3h	7Fh	0h
----	-----	----

HIT na via 0

ADDRESS (32 bits): 2FF0h

5h	7Fh	0h
----	-----	----

HIT na via 1

loop:

ADDRESS (32 bits): 1FECh

3h	7Eh	Ch
----	-----	----

MISS...

Preenche via 0

ADDRESS (32 bits): 2FECh

5h	7Eh	Ch
----	-----	----

MISS...

Preenche via 1

fim:

VA: .word <list of values>
VB: .word <list of values>

```
li    x10,1024 #N
la    x11,VA
la    x12,VB
addi  x13,x10,-1
sll   x13,x13,2
add   x14,x11,x13
add   x15,x12,x13
mv    x16,x0
ble   x10,x0,fim
lw    x17,0(x14)
lw    x18,0(x15)
mul   x17,x17,x18
add   x16,x16,x17
addi  x14,x14,-4
addi  x15,x15,-4
addi  x10,x10,-1
bgt   x10,x0,loop
```

Estrutura da cache

		TAG	V	DATA															
VIA 0			0																
			0																
			0																
			0																
	7Eh:	3h	1	VA[1019]				VA[1018]				VA[1017]				VA[1016]			
	7Fh:	3h	1	VA[1023]				VA[1022]				VA[1021]				VA[1020]			
				15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

		TAG	V	DATA															
VIA 1			0																
			0																
			0																
			0																
	7Eh:	5h	1	VB[1019]				VB[1018]				VB[1017]				VB[1016]			
	7Fh:	5h	1	VB[1023]				VB[1022]				VB[1021]				VB[1020]			
				15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- 2 vias (i.e., 2 tabelas)
- Capacidade para 4KB
- Linhas com 16B

Resultado:

Em cada 8 acessos:

- 4 para o array VA
- 4 para o array VB

Ocorrem:

- 2 **MISS** (um para cada array)
- 6 **HIT** (3 para cada array).

Assim:

$$HIT\ RATE = \frac{6}{8} = 75\%$$

```
VA:      .word  <list of values>
VB:      .word  <list of values>
```

```
li      x10,1024 #N
la      x11,VA
la      x12,VB
addi    x13,x10,-1
sll     x13,x13,2
add     x14,x11,x13
add     x15,x12,x13
mv      x16,x0
ble     x10,x0,fim
loop:   lw      x17,0(x14)
        lw      x18,0(x15)
        mul     x17,x17,x18
        add     x16,x16,x17
        addi    x14,x14,-4
        addi    x15,x15,-4
        addi    x10,x10,-1
        bgt     x10,x0,loop
```

fim:

Caches de mapeamento associativo

Vantagens e desvantagens do aumento da associatividade

- Em geral, quanto maior a associatividade, menor é o número de conflitos na cache

Embora na maioria dos casos esta regra seja verdade, é fácil arranjar exemplos onde uma cache com menos vias (ex: mapeamento direto) tem um maior HIT RATE que uma cache com mais vias (ex: completamente associativa)

- Contudo, o aumento do número de vias geralmente leva a:
 - Maior tempo de acesso (i.e., uma latência maior)
 - Mais recursos de hardware (é necessário lógica para leitura, comparação e seleção dos dados para cada uma das vias)
 - Maior área de silício e maior consumo de potência.

- Geralmente as caches têm um número de vias limitado (ex: 2-8 vias)
 - As caches completamente associativas, geralmente tem uma dimensão reduzida (ex: 16 ou 32 entradas)

Caches de mapeamento associativo (#vias>1)

Política de substituição

- Em caso de MISS, qual das vias deve ser escolhida?
 - Se a linha selecionada estiver livre ($V=0$) em pelo menos uma via → escolher essa via!
- O que fazer se essa linha já estiver ocupada em todas as vias?

Caches de mapeamento associativo (#vias>1)

Política de substituição

- Existem várias políticas de substituição:

- First-In First-Out (FIFO)**

➔ Substituir o bloco que está na cache há mais tempo

Requer a introdução de bits adicionais de controlo com a ordem de entrada na cache. Por exemplo, para uma cache com 4 vias, podemos usar 2 bits para indicar:

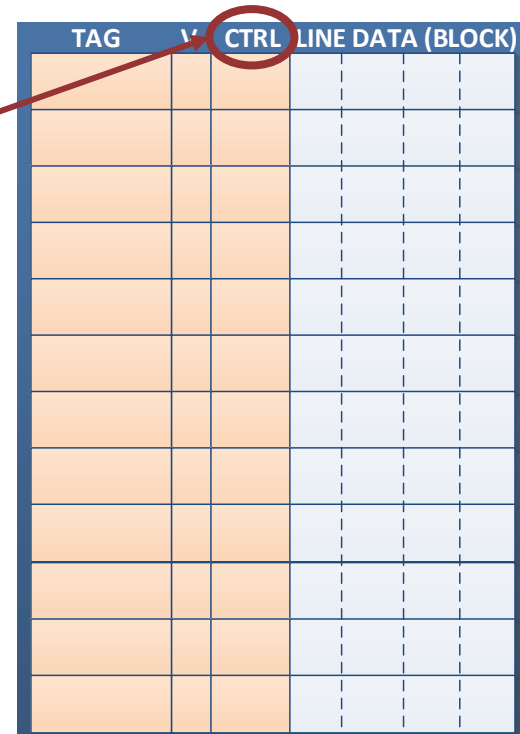
00 - inserido na cache há menos tempo

01

10

11 - inserido na cache há mais tempo (linha a substituir)

Sempre que acedemos a um bloco de uma via, teremos de alterar em todas as vias (i.e., vias 0,1,2,3), a ordem de entrada na cache



TAG	V	CTRL	LINE DATA (BLOCK)			

Caches de mapeamento associativo (#vias>1)

Política de substituição

- Existem várias políticas de substituição:

- Least Recently Used (LRU)**

➔ Substituir o bloco que está há mais tempo sem ser acedido

Requer a introdução de bits adicionais de controlo com a ordem de acesso. Por exemplo, para uma cache com 4 vias, podemos usar 2 bits para indicar:

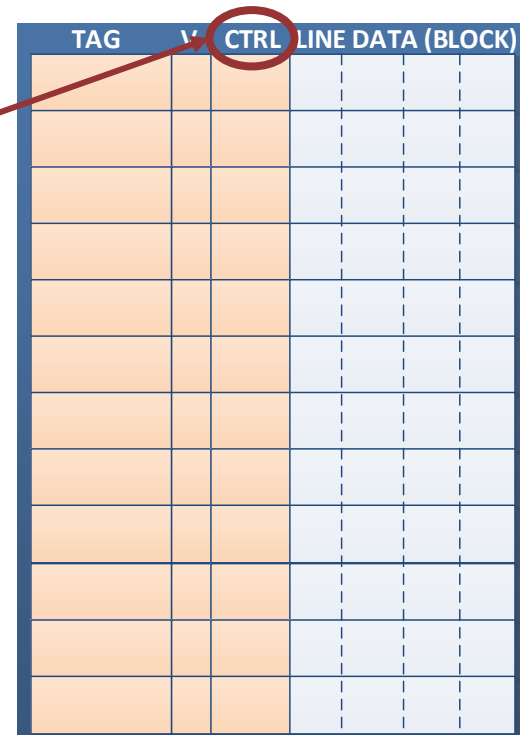
00 – usado há menos tempo

01

10

11 – usado há mais tempo (linha a substituir)

Sempre que acedemos a um bloco de uma via, teremos de alterar em todas as vias (i.e., vias 0,1,2,3), a ordem de acesso



TAG	V	CTRL	LINE DATA (BLOCK)			

Caches de mapeamento associativo (#vias>1)

Política de substituição

- Existem várias políticas de substituição:

- Random replacement**

- Escolher uma via ao acaso

Teoricamente simples e em geral não é significativamente pior que os outros dois.

Não requer manter qualquer estado adicional (i.e., não obriga há existência de mais bits de controlo)

TAG	V	CTRL	LINE DATA (BLOCK)			

Caches de mapeamento associativo (#vias>1)

Política de substituição

- Existem várias políticas de substituição:

- Pseudo-Least Recently Used (Pseudo-LRU)**

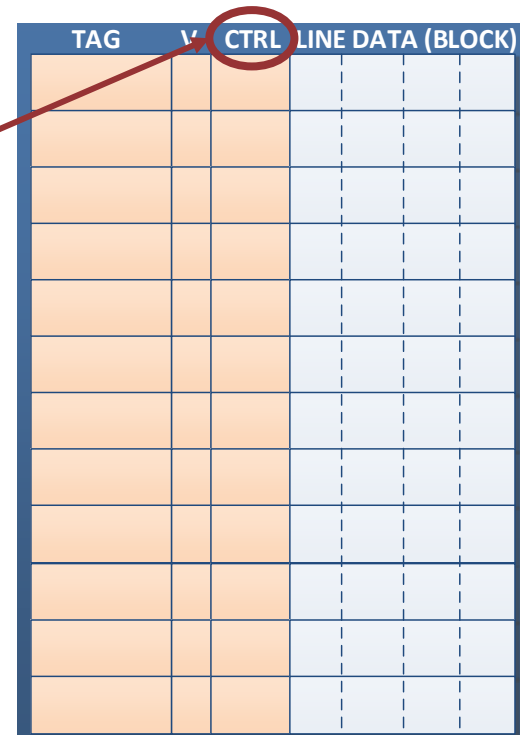
→ **Não** substituir o bloco que foi acedido há menos tempo

Requer a introdução um único bit adicional para indicar qual a via que contém o bloco acedido há menos tempo:

1 - último a ser usado

0- outros

Para #vias>2, esta política apenas indica qual a via cuja linha NÃO deve ser substituída. Entre as restantes vias alternativas, escolhemos uma política secundária (ex: random, round-robin)



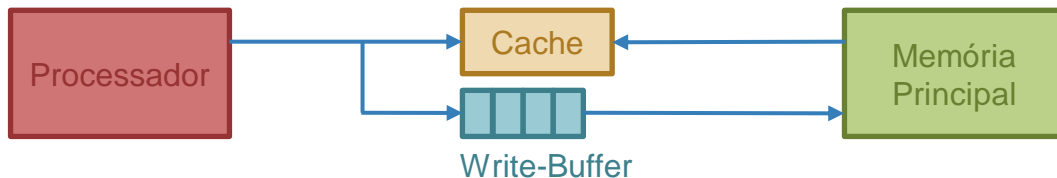
TAG	V	CTRL	LINE	DATA (BLOCK)

- Políticas de escrita:
 - **Write-back**: a escrita é realizada na cache
 - **Write-through**: a escrita não é realizada na cache.
(mas se o dado estiver em cache, este também é atualizado)

- Políticas de alocação:
 - **Write allocate**: uma escrita obriga à alocação dos dados na cache
→ conta para o número de acessos e gera um HIT/MISS
 - **Write not-allocate**: uma escrita nunca leva à alocação dos dados na cache
→ não conta para a contagem do número de acessos, e não gera HIT/MISS

- Combinações típicas:
 - Write-back, write allocate
 - Write-through, write not-allocate

- Políticas de escrita:
 - **Write-back**: a escrita é realizada na cache
 - **Write-through**: a escrita não é realizada na cache.
(mas se o dado estiver em cache, este também é atualizado)
- Para reduzir a latência das escritas, é habitual o uso de um write-buffer.



- O write-buffer permite que a cache continue a funcionar enquanto a escrita (demorada...) no nível seguinte da hierarquia de memória é realizado.
- Desta forma, as escritas passam a ser vistas como operações com latência nula.
- Problema: o que acontece se surgir uma instrução de LOAD logo a seguir ao STORE enquanto a escrita ainda está a decorrer?

Escrita de dados

Política de escrita e de alocação

- Políticas de escrita:
 - **Write-back:** a escrita é realizada na cache, mas não na memória principal
 - Ex: sw x5,0(x10)

CACHE

TAG	V	CTRL	LINE	DATA (BLOCK)
0134h	1			27

A2h

MEMÓRIA

0134A24h 13

Valor anterior

Escrita de dados

Política de escrita e de alocação

- Políticas de escrita:

- **Write-back:** a escrita é realizada na cache, mas não na memória principal

- Ex: sw x5,0(x10)

CACHE

Nova linha

00BC A2 0h



A2h

TAG	V	CTRL	LINE DATA (BLOCK)
0134h	1		27


- Substituição:

- Na sequência de um MISS na cache, a linha modificada anteriormente é substituída por uma outra linha
- Para garantir a coerência, é necessário escrever o valor na memória antes de alterar a linha!

- Contudo, é necessário que o controlador da cache tenha informação relativamente às linhas modificadas!
- **Solução:** Introduzir um bit de controlo adicional, **Dirty Bit (D)**, que indica quais as linhas modificadas.

MEMÓRIA

0134A24h



27

Escrita de dados

Política de escrita e de alocação

- Políticas de escrita:

- **Write-back:** a escrita é realizada na cache, mas não na memória principal

- Ex: sw x5,0(x10)

- **Dirty Bit (D):**

- O Dirty Bit (D=1) indica quais as linhas que devem ser escritas em memória antes de serem substituídas
- O Dirty Bit só é necessário para uma política de escrita do tipo Write-Back. Com uma política do tipo Write-Through, a memória (ou cache seguinte) já foi atualizada, assegurando desde logo a coerência dos dados.

CACHE

TAG	V	CTRL	LINE DATA (BLOCK)			
0134h	1					27

Nova linha

00BC A2 0h



A2h

MEMÓRIA

0134A24h	27



Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

EXERCÍCIO #7

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Admita os seguintes endereços para a localização dos vetores VA, VB e VC:

VA = 1000h; VB = 2000h; VC = 3000h

e ainda as seguintes latências de acessos aos dados:

Latência da L1 = 3 ciclos

Latência da DRAM = 100 ciclos

1. **Determine qual a combinação de política de escrita e alocação de dados mais vantajosa.**
2. **Para cada caso, determine o número ótimo de vias, de forma a minimizar o número de MISSes.**

```
float *A, *B, *C;
int i;
...
for (i=0; i<1024; i++)
    C[i] = B[i]*(i>0?B[i-1]:1.0) + A[i];
```

A: .float <list of values>
B: .float <list of values>
C: .zero 4*1024

```
li      x10,1024    #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0       #i=0
addi    x14,x0,1
fcvt.s.w f0,x14     #f0=1.0
bgt     x5,x10,end  #i>N?
loop:   flw         f1,0(x11)  #A[i]
        flw         f2,0(x12)  #B[i]
        fmv.s      f3,f0
        blt        x5,x0,skip  #i<=0?
        flw         f3,-4(x12)
skip:   fmul.s      f2,f2,f3
        fadd.s     f1,f2,f1
        fsw        f1,0(x13)   #C[i]=
        addi       x11,x11,4
        addi       x12,x12,4
        addi       x13,x13,4
        addi       x5,x5,1     #i++
        blt        x5,x10,loop #i<N?

end:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Linhas com 16B → offset de 4 bits

Capacidade para 4 KB → 256 linhas

Número máximo de bits para índice: 8 (cache de mapeamento direto)

```
A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero       4*1024

li      x10,1024    #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0       #i=0
addi    x14,x0,1
fcvt.s.w f0,x14     #f0=1.0
bgt     x5,x10,end  #i>N?
loop:   flw         f1,0(x11)  #A[i]
        flw         f2,0(x12)  #B[i]
        fmv.s       f3,f0
        blt         x5,x0,skip  #i<=0?
        flw         f3,-4(x12)
skip:   fmul.s      f2,f2,f3
        fadd.s      f1,f2,f1
        fsw         f1,0(x13)   #C[i]=
        addi        x11,x11,4
        addi        x12,x12,4
        addi        x13,x13,4
        addi        x5,x5,1     #i++
        blt         x5,x10,loop #i<N?

end:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Back, Write Allocate

Sequência de acessos:

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	1000	2000	-	3000
1	1004	2004	2000	3004
2	1008	2008	2004	3008
3	100C	200C	2008	300C
4	1010	2010	200C	3010
5	1014	2014	2010	3014
6	1018	2018	2014	3018
7	101C	201C	2018	301C
8	1020	2020	201C	3020
9	1024	2024	2000	3024
10	1028	2028	2004	3028
11	102C	202C	2008	302C
...

```
float *A, *B, *C;
int i;
...
for (i=0; i<1024; i++)
    C[i] = B[i]*(i>0?B[i-1]:1.0) + A[i];
```

A: .float <list of values>
B: .float <list of values>
C: .zero 4*1024

```
li x10,1024 #N
la x11,VA
la x12,VB
la x13,VC
mv x5,x0 #i=0
addi x14,x0,1
fcvt.s.w f0,x14 #f0=1.0
bgt x5,x10,end #i>N?
loop: flw f1,0(x11) #A[i]
      flw f2,0(x12) #B[i]
      fmv.s f3,f0
      blt x5,x0,skip #i<=0?
      flw f3,-4(x12)
skip: fmul.s f2,f2,f3
      fadd.s f1,f2,f1
      fsw f1,0(x13) #C[i]=
      addi x11,x11,4
      addi x12,x12,4
      addi x13,x13,4
      addi x5,x5,1 #i++
      blt x5,x10,loop #i<N?
end:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Back, Write Allocate

Sequência de acessos:

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	100 0	200 0	-	300 0
1	100 4	200 4	200 0	300 4
2	100 8	200 8	200 4	300 8
3	100 C	200 C	200 8	300 C
4	101 0	201 0	200 C	301 0
5	101 4	201 4	201 0	301 4
6	101 8	201 8	201 4	301 8
7	101 C	201 C	201 8	301 C
8	102 0	202 0	201 C	302 0
9	102 4	202 4	200 0	302 4
10	102 8	202 8	200 4	302 8
11	102 C	202 C	200 8	302 C
...

Para uma cache de mapeamento direto, todos os acessos numa dada iteração acedem à mesma linha...

```

A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero        4*1024

li      x10,1024      #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0          #i=0
addi    x14,x0,1
fcvt.s.w f0,x14        #f0=1.0
bgt     x5,x10,end     #i>N?
loop:   flw          f1,0(x11)  #A[i]
        flw          f2,0(x12)  #B[i]
        fmv.s        f3,f0
        blt          x5,x0,skip  #i<=0?
        flw          f3,-4(x12)
skip:   fmul.s        f2,f2,f3
        fadd.s        f1,f2,f1
        fsw          f1,0(x13)  #C[i]=
        addi          x11,x11,4
        addi          x12,x12,4
        addi          x13,x13,4
        addi          x5,x5,1    #i++
        blt          x5,x10,loop #i<N?

end:

```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Back, Write Allocate

Sequência de acessos:

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	100 0	200 0	-	300 0
1	100 4	200 4	200 0	300 4
2	100 8	200 8	200 4	300 8
3	100 C	200 C	200 8	300 C
4	101 0	201 0	200 C	301 0
5	101 4	201 4	201 0	301 4
6	101 8	201 8	201 4	301 8
7	101 C	201 C	201 8	301 C
8	102 0	202 0	201 C	302 0
9	102 4	202 4	200 0	302 4
10	102 8	202 8	200 4	302 8
11	102 C	202 C	200 8	302 C
...

Para evitar conflitos, são precisas pelo menos 3 vias.
Como 256 linhas não podem ser distribuídas em 3 vias, são precisas pelo menos 4 vias.

```

A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero        4*1024

li      x10,1024      #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0         #i=0
addi    x14,x0,1
fcvt.s.w f0,x14       #f0=1.0
bgt     x5,x10,end    #i>N?
loop:   flw          f1,0(x11)  #A[i]
        flw          f2,0(x12)  #B[i]
        fmv.s        f3,f0
        blt          x5,x0,skip  #i<=0?
        flw          f3,-4(x12)
skip:   fmul.s        f2,f2,f3
        fadd.s        f1,f2,f1
        fsw          f1,0(x13)  #C[i]=
        addi         x11,x11,4
        addi         x12,x12,4
        addi         x13,x13,4
        addi         x5,x5,1     #i++
        blt          x5,x10,loop #i<N?

end:

```


Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Back, Write Allocate

Resultado dos acessos

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	1 00 0 M	2 00 0 M	-	3 00 0 M
1	1 00 4 H	2 00 4 H	2 00 0 H	3 00 4 H
2	1 00 8 H	2 00 8 H	2 00 4 H	3 00 8 H
3	1 00 C H	2 00 C H	2 00 8 H	3 00 C H
4	1 01 0 M	2 01 0 M	2 00 C H	3 01 0 M
5	1 01 4 H	2 01 4 H	2 01 0 H	3 01 4 H
6	1 01 8 H	2 01 8 H	2 01 4 H	3 01 8 H
7	1 01 C H	2 01 C H	2 01 8 H	3 01 C H
8	1 02 0 M	2 02 0 M	2 01 C H	3 02 0 M
9	1 02 4 H	2 02 4 H	2 00 0 H	3 02 4 H
10	1 02 8 H	2 02 8 H	2 00 4 H	3 02 8 H
11	1 02 C H	2 02 C H	2 00 8 H	3 02 C H
...

M MISS

H HIT

A: .float <list of values>
B: .float <list of values>
C: .zero 4*1024

```
li      x10,1024    #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0       #i=0
addi    x14,x0,1
fcvt.s.w f0,x14     #f0=1.0
bgt     x5,x10,end  #i>N?
loop:   flw         f1,0(x11)  #A[i]
        flw         f2,0(x12)  #B[i]
        fmv.s       f3,f0
        blt         x5,x0,skip #i<=0?
        flw         f3,-4(x12)
skip:   fmul.s      f2,f2,f3
        fadd.s      f1,f2,f1
        fsw         f1,0(x13)  #C[i]=
        addi        x11,x11,4
        addi        x12,x12,4
        addi        x13,x13,4
        addi        x5,x5,1    #i++
        blt         x5,x10,loop #i<N?
```

end:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Back, Write Allocate

MISS RATE = $768/4092=18.77\%$

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	1 00 0 M	2 00 0 M	-	3 00 0 M
1	1 00 4 H	2 00 4 H	2 00 0 H	3 00 4 H
2	1 00 8 H	2 00 8 H	2 00 4 H	3 00 8 H
3	1 00 C H	2 00 C H	2 00 8 H	3 00 C H
4	1 01 0 M	2 01 0 M	2 00 C H	3 01 0 M
5	1 01 4 H	2 01 4 H	2 01 0 H	3 01 4 H
6	1 01 8 H	2 01 8 H	2 01 4 H	3 01 8 H
7	1 01 C H	2 01 C H	2 01 8 H	3 01 C H
8	1 02 0 M	2 02 0 M	2 01 C H	3 02 0 M
9	1 02 4 H	2 02 4 H	2 00 0 H	3 02 4 H
10	1 02 8 H	2 02 8 H	2 00 4 H	3 02 8 H
11	1 02 C H	2 02 C H	2 00 8 H	3 02 C H
...

Loop: 256 grupos de 4 iterações

Grupo #1 (i=0,1,2,3): 3 MISSes, 15 acessos

Grupo #n (i>3): 3 MISSes, 16 acessos

```

A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero       4*1024
  
```

```

li      x10,1024    #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0       #i=0
addi    x14,x0,1
fcvt.s.w f0,x14     #f0=1.0
bgt     x5,x10,end  #i>N?
loop:   flw         f1,0(x11)  #A[i]
        flw         f2,0(x12)  #B[i]
        fmv.s       f3,f0
        blt         x5,x0,skip #i<=0?
        flw         f3,-4(x12)
skip:   fmul.s      f2,f2,f3
        fadd.s      f1,f2,f1
        fsw         f1,0(x13)  #C[i]=
        addi        x11,x11,4
        addi        x12,x12,4
        addi        x13,x13,4
        addi        x5,x5,1    #i++
        blt         x5,x10,loop #i<N?
  
```

end:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Back, Write Allocate

MISS RATE = $768/4092=18.77\%$

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	1 00 0 M	2 00 0 M	-	3 00 0 M
1	1 00 4 H	2 00 4 H	2 00 0 H	3 00 4 H
2	1 00 8 H	2 00 8 H	2 00 4 H	3 00 8 H
3	1 00 C H	2 00 C H	2 00 8 H	3 00 C H
4	1 01 0 M	2 01 0 M	2 00 C H	3 01 0 M
5	1 01 4 H	2 01 4 H	2 01 0 H	3 01 4 H
6	1 01 8 H	2 01 8 H	2 01 4 H	3 01 8 H
7	1 01 C H	2 01 C H	2 01 8 H	3 01 C H
8	1 02 0 M	2 02 0 M	2 01 C H	3 02 0 M
9	1 02 4 H	2 02 4 H	2 00 0 H	3 02 4 H
10	1 02 8 H	2 02 8 H	2 00 4 H	3 02 8 H
11	1 02 C H	2 02 C H	2 00 8 H	3 02 C H
...

Loop: 256 grupos de 4 iterações

Grupo #1 (i=0,1,2,3): 3 MISSes, 15 acessos

Grupo #n (i>3): 3 MISSes, 16 acessos

Miss Rate = $(3+3*255)/(15+16*255)=18.71\%$

```

A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero       4*1024

```

```

li      x10,1024    #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0       #i=0
addi    x14,x0,1
fcvt.s.w f0,x14     #f0=1.0
bgt     x5,x10,end  #i>N?
loop:   flw         f1,0(x11)  #A[i]
        flw         f2,0(x12)  #B[i]
        fmv.s      f3,f0
        blt        x5,x0,skip  #i<=0?
        flw         f3,-4(x12)
skip:   fmul.s      f2,f2,f3
        fadd.s     f1,f2,f1
        fsw        f1,0(x13)   #C[i]=
        addi       x11,x11,4
        addi       x12,x12,4
        addi       x13,x13,4
        addi       x5,x5,1     #i++
        blt        x5,x10,loop #i<N?

```

end:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Through, Write Not-Allocate

Sequência de acessos:

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	1000	2000	-	-
1	1004	2004	2000	-
2	1008	2008	2004	-
3	100C	200C	2008	-
4	1010	2010	200C	-
5	1014	2014	2010	-
6	1018	2018	2014	-
7	101C	201C	2018	-
8	1020	2020	201C	-
9	1024	2024	2000	-
10	1028	2028	2004	-
11	102C	202C	2008	-
...	-

Como os stores não contam (devido à presença do **write-buffer**), bastam 2 vias

```

A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero        4*1024

li      x10,1024      #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0         #i=0
addi    x14,x0,1
fcvt.s.w f0,x14       #f0=1.0
bgt     x5,x10,end    #i>N?
loop:   flw          f1,0(x11)  #A[i]
        flw          f2,0(x12)  #B[i]
        fmv.s        f3,f0
        blt          x5,x0,skip  #i<=0?
        flw          f3,-4(x12)
skip:   fmul.s        f2,f2,f3
        fadd.s        f1,f2,f1
        fsw          f1,0(x13)  #C[i]=
        addi         x11,x11,4
        addi         x12,x12,4
        addi         x13,x13,4
        addi         x5,x5,1     #i++
        blt          x5,x10,loop #i<N?

end:

```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Through, Write Not-Allocate

Resultado dos acessos

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	1000 M	2000 M	-	-
1	1004 H	2004 H	2000 H	-
2	1008 H	2008 H	2004 H	-
3	100C H	200C H	2008 H	-
4	1010 M	2010 M	200C H	-
5	1014 H	2014 H	2010 H	-
6	1018 H	2018 H	2014 H	-
7	101C H	201C H	2018 H	-
8	1020 M	2020 M	201C H	-
9	1024 H	2024 H	2000 H	-
10	1028 H	2028 H	2004 H	-
11	102C H	202C H	2008 H	-
...	-

M MISS

H HIT

```
A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero        4*1024
```

```
li      x10,1024      #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0          #i=0
addi    x14,x0,1
fcvt.s.w f0,x14        #f0=1.0
bgt     x5,x10,end     #i>N?
loop:   flw            f1,0(x11)  #A[i]
        flw            f2,0(x12)  #B[i]
        fmv.s          f3,f0
        blt            x5,x0,skip #i<=0?
        flw            f3,-4(x12)
skip:   fmul.s         f2,f2,f3
        fadd.s          f1,f2,f1
        fsw             f1,0(x13)  #C[i]=
        addi            x11,x11,4
        addi            x12,x12,4
        addi            x13,x13,4
        addi            x5,x5,1    #i++
        blt            x5,x10,loop #i<N?
```

end:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Through, Write Not-Allocate

Resultado dos acessos

i	lw A[i]	lw B[i]	lw B[i-1]	sw C[i]
0	1 00 0 M	2 00 0 M	-	-
1	1 00 4 H	2 00 4 H	2 00 0 H	-
2	1 00 8 H	2 00 8 H	2 00 4 H	-
3	1 00 C H	2 00 C H	2 00 8 H	-
4	1 01 0 M	2 01 0 M	2 00 C H	-
5	1 01 4 H	2 01 4 H	2 01 0 H	-
6	1 01 8 H	2 01 8 H	2 01 4 H	-
7	1 01 C H	2 01 C H	2 01 8 H	-
8	1 02 0 M	2 02 0 M	2 01 C H	-
9	1 02 4 H	2 02 4 H	2 00 0 H	-
10	1 02 8 H	2 02 8 H	2 00 4 H	-
11	1 02 C H	2 02 C H	2 00 8 H	-
...	-

Loop: 256 grupos de 4 iterações

Grupo #1 (i=0,1,2,3): 2 MISSes, 11 acessos

Grupo #n: 2 MISSes, 12 acessos

Miss Rate = $(2+2*255)/(11+12*255) = 16.67\%$

```
A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero        4*1024
```

```
li      x10,1024    #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0        #i=0
addi    x14,x0,1
fcvt.s.w f0,x14      #f0=1.0
bgt     x5,x10,end   #i>N?
loop:   flw          f1,0(x11)  #A[i]
        flw          f2,0(x12)  #B[i]
        fmv.s        f3,f0
        blt          x5,x0,skip #i<=0?
        flw          f3,-4(x12)
skip:   fmul.s       f2,f2,f3
        fadd.s        f1,f2,f1
        fsw          f1,0(x13)  #C[i]=
        addi          x11,x11,4
        addi          x12,x12,4
        addi          x13,x13,4
        addi          x5,x5,1    #i++
        blt          x5,x10,loop #i<N?
```

end:

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 4KB
- Linhas com 16B

Write-Back, Write Allocate

MISS RATE $\approx 18.71\%$

Tempo médio de acesso aos dados:

$$T = 3 + 0.1871 \times 100 = 21.71 \text{ ciclos}$$

Latência total de acesso (~#stalls):

Assumindo que sempre que temos miss, vamos ter *stalls* até que o pedido seja satisfeito, mas que os *hits* não geram *stalls*.

$$\text{Lat}_{\text{Total}} \approx 18.71 \text{ stalls/acesso} \\ \times 4095 \text{ acessos} \approx 77\text{k ciclos}$$

Write-Through, Write Not-Allocate

MISS RATE $\approx 16.67\%$

Tempo médio de acesso aos dados:

$$T = 3 + 0.1667 \times 100 = 19.67 \text{ ciclos}$$

Latência total de acesso (~#stalls):

$$\text{Lat}_{\text{Total}} = 16.67 \times 3071 \approx 51\text{k ciclos}$$

```
A:      .float      <list of values>
B:      .float      <list of values>
C:      .zero        4*1024
```

```
li      x10,1024      #N
la      x11,VA
la      x12,VB
la      x13,VC
mv      x5,x0          #i=0
addi    x14,x0,1
fcvt.s.w f0,x14        #f0=1.0
bgt     x5,x10,end     #i>N?
loop:   flw            f1,0(x11)  #A[i]
        flw            f2,0(x12)  #B[i]
        fmv.s          f3,f0
        blt            x5,x0,skip #i<=0?
        flw            f3,-4(x12)
skip:   fmul.s          f2,f2,f3
        fadd.s          f1,f2,f1
        fsw             f1,0(x13)  #C[i]=
        addi            x11,x11,4
        addi            x12,x12,4
        addi            x13,x13,4
        addi            x5,x5,1    #i++
        blt            x5,x10,loop #i<N?

end:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 16KB
- Linhas com 32B

EXERCÍCIO #8

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variáveis alocadas em registos
```

```
for ( i=0 ; i<2048 ; i+=2 )
    C[i] = 2*A[i] + 4*B[i];
```

```
for ( i=0 ; i<1024 ; i+=1 )
    C[2*i] = C[2*i] * D[2*i];
```

- Esboce a organização da memória cache, indicando a decomposição dos bits de endereço em TAG, INDEX e OFFSET
- Indique a taxa de falhas para o troço de código indicado, assumindo que as variáveis são lidas da memória e escritas em memória pela ordem indicada no código.

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variáveis alocadas em registos
```

```
for ( i=0 ; i<2048 ; i+=2 )
```

```
    C[i] = 2*A[i] + 4*B[i];
```

3

1

2

A[0],B[0],C[0],A[2],B[2],C[2],...

```
for ( i=0 ; i<1024 ; i+=1 )
```

```
    C[2*i] = C[2*i] * D[2*i];
```

3

1

2

C[0],D[0],C[0],C[2],D[2],C[2],...

- Esboce a organização da memória cache, indicando a decomposição dos bits de endereço em TAG, INDEX e OFFSET
- Indique a taxa de falhas para o troço de código indicado, **assumindo que as variáveis são lidas da memória e escritas em memória pela ordem indicada no código.**

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variáveis alocadas em registos
```

```
for ( i=0 ; i<2048 ; i+=2 )
    C[i] = 2*A[i] + 4*B[i];
```

```
for ( i=0 ; i<1024 ; i+=1 )
    C[2*i] = C[2*i] * D[2*i];
```

- a) Esboce a organização da memória cache, indicando a decomposição dos bits de endereço em TAG, INDEX e OFFSET

TAG=20 bits, INDEX=6 bits, OFFSET=6 bits

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

Address

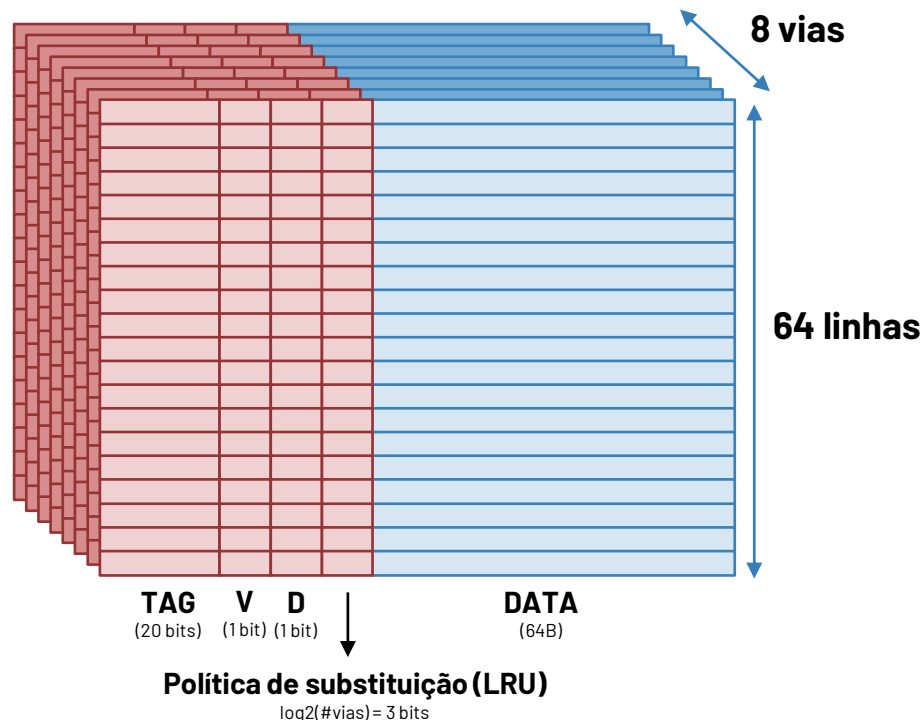
TAG:
20 bits

INDEX:
6 bits

OFFSET:
6 bits

Estrutura da cache:

- 8 vias (tabelas) todas com a mesma estrutura.



Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

Address

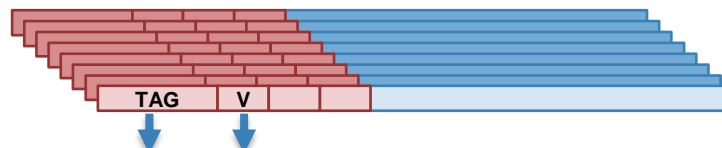
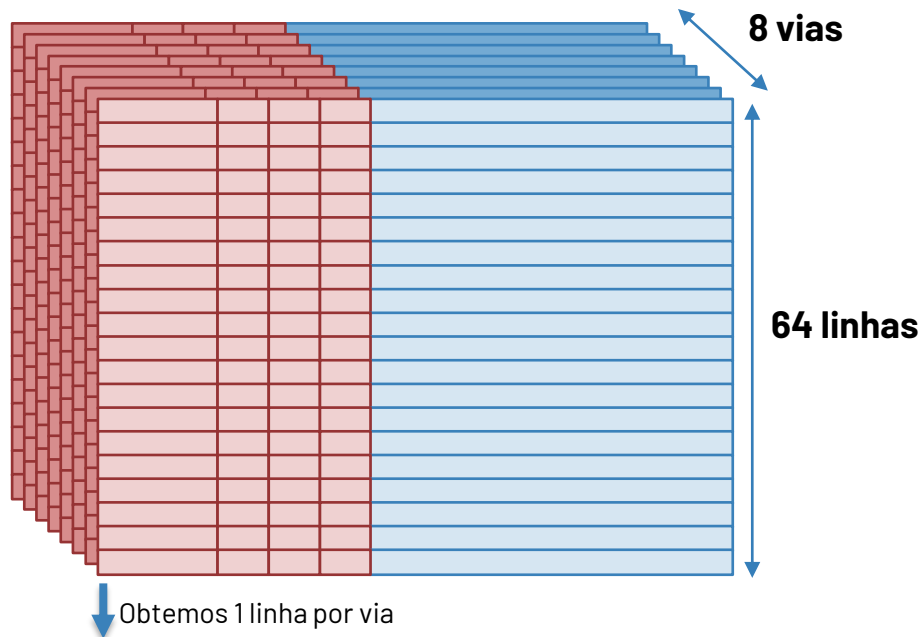
TAG:	INDEX:	OFFSET:
20 bits	6 bits	6 bits

Index:

Escolhe a linha a ler em cada via (tabela)

Estrutura da cache:

- 8 vias (tabelas) todas com a mesma estrutura.



Comparamos as TAGs de cada linha e verificamos os bits de validade. Usamos tantos comparadores, quanto o número de vias.

Se houver hit numa das VIAS escolhemos essa via. Caso contrário, vamos pedir a linha correspondente ao nível seguinte (cache/memória)

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

Address

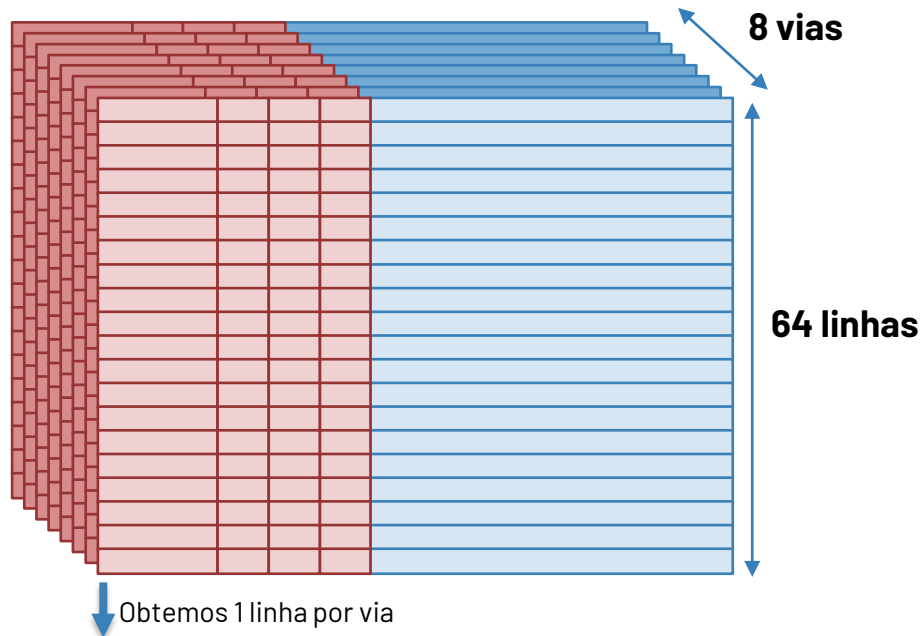
TAG:	INDEX:	OFFSET:
20 bits	6 bits	6 bits

Index:

Escolhe a linha a ler em cada via (tabela)

Estrutura da cache:

- 8 vias (tabelas) todas com a mesma estrutura.



Offset: seleciona o primeiro byte (dentro do bloco de dados) da palavra a ler

Via com HIT

TAG	V		
-----	---	--	--

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variáveis alocadas em registos
```

```
for ( i=0 ; i<2048 ; i+=2 )
    C[i] = 2*A[i] + 4*B[i];
```

```
for ( i=0 ; i<1024 ; i
    C[2*i] = C[2*i] *
```

Começa por aceder a A[0]: linha 0, offset 0. Gera um MISS. Vamos preencher por exemplo a Via 0. Como a linha tem capacidade para 64B, vamos preencher a linha 0 com 16 palavras: A[0],A[1],...,A[15].

Como a variável i é incrementada em 2 em cada iteração do loop, apenas os elementos pares são úteis.

- b) Indique a taxa de falhas para o troço de código indicado, assumindo que as variáveis são lidas da memória e escritas em memória pela ordem indicada no código.

Assumindo que a cache está inicialmente vazia...

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variáveis alocadas em registos
```

```
for ( i=0 ; i<2048 ; i+=2 )
    C[i] = 2*A[i] + 4*B[i];
```

Acedemos a B[0]: linha 0, offset 0. Gera um MISS.

```
for ( i=0 ; i<1024 ; i+=1 )
    C[2*i] = C[2*i] * D[2*i];
```

Como a via 0 está ocupada, vamos preencher a via 1 com 16 palavras: B[0],B[1],...,B[15].

- b) Indique a taxa de falhas para o troço de código indicado, assumindo que as variáveis são lidas da memória e escritas em memória pela ordem indicada no código.

Assumindo que a cache está inicialmente vazia...

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variáveis alocadas em registos
```

```
for ( i=0 ; i<2048 ; i+=2 )
```

```
    C[i] = 2*A[i] + 4*B[i];
```

Acedemos a C[0]: linha 0, offset 0. Gera um MISS.

```
for ( i=0 ; i<1024 ; i+=1 )
```

```
    C[2*i] = C[2*i] * D[2*i]
```

Como as vias 0 e 1 estão ocupadas, vamos preencher a via 2 com 16 palavras: C[0],C[1],...,C[15]. De seguida atualizamos C[0] com o novo valor.

- b) Indique a taxa de falhas para o troço de código indicado, assumindo que as variáveis são lidas da memória e escritas em memória pela ordem indicada no código.

Assumindo que a cache está inicialmente vazia...

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variável de índice
```

```
for ( i=0 ; i<2048 ; i+=2 )
```

$C[i] = 2 \cdot A[i] + 4 \cdot B[i]$

```
for ( i=0 ; i<1024 ; i+=1 )
    C[2*i] = C[2*i] * D[2*i]
```

2ª a 7ª iteração do loop:

- A[2]....A[14]: linha 0, offset 0. 7 HITs na via 0
- B[2]....B[14]: linha 0, offset 0. 7 HITs na via 1
- C[2]....C[14]: linha 0, offset 0. 7 HITs na via 2

Resultado:

A cada 8 iterações temos:

- 3x1 MISS (uma falha por cada vector)
- 3x7 HITs

NOTA:

- Cada vetor tem 2048 elementos x 4 B = 8KB
- Cada via da cache tem 32KB/8 = 4KB.
- Na via 0 cabem A[0] a A[1023]. Quando chegamos a A[1024] precisamos de ocupar uma nova via.
- A[1024]...A[2047] → Via 3
- B[1024]...B[2047] → Via 4
- C[1024]...C[2047] → Via 5

- b) Indique a taxa de falhas para o troço as variáveis são lidas da memória e é indicada no código.

Assumindo que a cache está inicialmente vazia...

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variáveis alocadas em registos
```

```
for ( i=0 ; i<2048 ; i+=2 )
    C[i] = 2*A[i] + 4*B[i];
```

Quando chegamos ao segundo loop temos os vetores A, B e C na cache.

```
for ( i=0 ; i<1024 ; i+=1 )
    C[2*i] = C[2*i] * D[2*i];
```

O vetor D vai ocupar as vias 6 e 7. A cada miss temos 7 hits.

O acesso ao vetor C vai dar sempre HIT.

- b) Indique a taxa de falhas para o troço de código indicado, assumindo que as variáveis são lidas da memória e escritas em memória pela ordem indicada no código.

Assumindo que a cache está inicialmente vazia...

Mapeamento dos dados na cache

Cache L1 de dados:

- 8 vias
- dimensão de 32KB
- blocos de 64B
- Política de escrita: *write back, write allocate*
- Política de substituição: *Least Recently Used (LRU)*

```
int A[2048];           // endereço base: 0x0F 0000
int B[2048];           // endereço base: 0x0F 1000
int C[2048];           // endereço base: 0x0F 2000
int D[2048];           // endereço base: 0x0F 3000
register int i;         // variáveis alocadas em registos
```

```
for ( i=0 ; i<2048 ; i+=2 )
    C[i] = 2*A[i] + 4*B[i];
```

3x1024 acessos
3x1024/8 Misses

```
for ( i=0 ; i<1024 ; i+=1 )
    C[2*i] = C[2*i] * D[2*i];
```

3x1024 acessos
1024/8 Misses (todos no vector D)

- b) Indique a taxa de falhas para o troço de código indicado, assumindo que as variáveis são lidas da memória e escritas em memória pela ordem indicada no código.

Taxa de falhas:

$$MISS\ RATE = \frac{3 \times \frac{1024}{8} + \frac{1024}{8}}{3 \times 1024 + 3 \times 1024} = \frac{1}{12} \approx 8.33\%$$

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 32B
- Linhas com 8B
- Política de substituição:
Least Recently Used
(LRU)

EXERCÍCIO #9

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 32B
- Linhas com 8B
- Política de substituição: *Least Recently Used (LRU)*

Compare o desempenho de uma cache de mapeamento direto e uma cache completamente associativa.

Para tal, determine a taxa de sucesso no acesso às instruções para o troço de código indicado para cada um dos casos.

Indique ainda o conteúdo da cache após a primeira execução da instrução "B.GT loop"

A: .float <list of values>
B: .float <list of values>
C: .zero 4*1024

Address:

```
00h      li      x10,1024    #N
04h      la      x11,VA
08h      la      x12,VB
0Ch      la      x13,VC
10h      mv      x5,x0       #i=0
14h      addi    x14,x0,1
18h      fcvt.s.w f0,x14     #f0=1.0
1Ch      bgt     x5,x10,end  #i>N?
20h loop: flw     f1,0(x11)   #A[i]
24h      flw     f2,0(x12)   #B[i]
28h      fmv.s   f3,f0
2Ch      ble     x5,x0,skip  #i<=0?
30h      flw     f3,-4(x12)
34h skip: fmul.s  f2,f2,f3
38h      fadd.s  f1,f2,f1
3Ch      fsw     f1,0(x13)   #C[i]=
40h      addi    x11,x11,4
44h      addi    x12,x12,4
48h      addi    x13,x13,4
4Ch      addi    x5,x5,1     #i++
50h      blt     x5,x10,loop #i<N?
54h end:
```

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 32B → **4 linhas**
- Linhas com 8B → **2 inst./linha**
- Política de substituição: *Least Recently Used (LRU)*

PRIMEIRA ITERAÇÃO

Map. Direto			Completamente Associativa					
RES.	LINHA	TAG	RES.	VIA	TAG	Address:		
M	0	0	M	0	0h	00h	li	x10,1024 #N
H	0	0	H	0	0h	04h	la	x11,VA
M	1	0	M	1	1h	08h	la	x12,VB
H	1	0	H	1	1h	0Ch	la	x13,VC
M	2	0	M	2	2h	10h	mv	x5,x0 #i=0
H	2	0	H	2	2h	14h	addi	x14,x0,1
M	3	0	M	3	3h	18h	fcvt.s.w	f0,x14 #f0=1.0
H	3	0	H	3	3h	1Ch	bgt	x5,x10,end #i>N?
M	0	1	M	0	4h	20h	loop: flw	f1,0(x11) #A[i]
H	0	1	H	0	4h	24h	flw	f2,0(x12) #B[i]
M	1	1	M	1	5h	28h	fmv.s	f3,f0
H	1	1	H	1	5h	2Ch	blt	x5,x0,skip #i<=0?
						30h	flw	f3,-4(x12)
M	2	1	M	2	6h	34h	skip: fmul.s	f2,f2,f3
M	3	1	M	3	7h	38h	fadd.s	f1,f2,f1
H	3	1	H	3	7h	3Ch	fsw	f1,0(x13) #C[i]=
M	0	2	M	0	8h	40h	addi	x11,x11,4
H	0	2	H	0	8h	44h	addi	x12,x12,4
M	1	2	M	1	9h	48h	addi	x13,x13,4
H	1	2	H	1	9h	4Ch	addi	x5,x5,1 #i++
M	2	2	M	2	Ah	50h	blt	x5,x10,loop #i<N?
						54h	end:	

A: .float <list of values>
 B: .float <list of values>
 C: .zero 4*1024

Estrutura da cache

Após a 1ª iteração do loop:

Mapeamento direto

<u>Via única:</u>	TAG	V	DATA							
Linha 0:	2h	1	addi x12,x12,4				addi x11,x11,4			
Linha 1:	2h	1	addi x5,x5,1				addi x13,x13,4			
Linha 2:	2h	1	???				blt x5,x10,loop			
Linha 3:	1h	1	fsw f1,0(x13)				fadd.s f1,f2,f1			
			7	6	5	4	3	2	1	0

Hits: 9
Misses: 11
Acessos: 20

Completamente associativa

<u>1 linha/via</u>	TAG	V	DATA							
Via 0:	8h	1	addi x12,x12,4				addi x11,x11,4			
Via 1:	9h	1	addi x5,x5,1				addi x13,x13,4			
Via 2:	Ah	1	???				blt x5,x10,loop			
Via 3:	7h	1	fsw f1,0(x13)				fadd.s f1,f2,f1			
			7	6	5	4	3	2	1	0

Hits: 9
Misses: 11
Acessos: 20

Mapeamento dos dados na cache

Considere uma cache L1 para um processador com ISA RV32G

Características da memória cache:

- Capacidade para 32B → **4 linhas**
- Linhas com 8B → **2 inst./linha**
- Política de substituição: *Least Recently Used (LRU)*

SEGUNDA ITERAÇÃO

Map. Direto			Completamente Associativa					
RES.	LINHA	TAG	RES.	VIA	TAG	Address:		
						00h	li	x10,1024 #N
						04h	la	x11,VA
						08h	la	x12,VB
						0Ch	la	x13,VC
						10h	mv	x5,x0 #i=0
						14h	addi	x14,x0,1
						18h	fcvt.s.w	f0,x14 #f0=1.0
						1Ch	bgt	x5,x10,end #i>N?
M	0	1	M	3	4h	20h	loop: flw	f1,0(x11) #A[i]
H	0	1	H	3	4h	24h	flw	f2,0(x12) #B[i]
M	1	1	M	0	5h	28h	fmv.s	f3,f0
H	1	1	H	0	5h	2Ch	blt	x5,x0,skip #i<=0?
M	2	1	M	1	6h	30h	flw	f3,-4(x12)
H	2	1	H	1	6h	34h	skip: fmul.s	f2,f2,f3
H	3	1	M	2	7h	38h	fadd.s	f1,f2,f1
H	3	1	H	2	7h	3Ch	fsw	f1,0(x13) #C[i]=
M	0	2	M	3	8h	40h	addi	x11,x11,4
H	0	2	H	3	8h	44h	addi	x12,x12,4
M	1	2	M	0	9h	48h	addi	x13,x13,4
H	1	2	H	0	9h	4Ch	addi	x5,x5,1 #i++
M	2	2	M	1	Ah	50h	blt	x5,x10,loop #i<N?
						54h	end:	

A: .float <list of values>
 B: .float <list of values>
 C: .zero 4*1024

Estrutura da cache

Após a 2ª iteração do loop:

Mapeamento direto

Via única:	TAG	V	DATA								
	0:	2h	1	addi x12,x12,4				addi x11,x11,4			
	1:	2h	1	addi x5,x5,1				addi x13,x13,4			
	2:	2h	1	???				blt x5,x10,loop			
	3:	1h	1	fsw f1,0(x13)				fadd.s f1,f2,f1			
				7	6	5	4	3	2	1	0

Hits: 16
Misses: 17
Acessos: 33

Completamente associativa

	TAG	V	DATA							
Via 0:	9h	1	addi x5,x5,1				addi x13,x13,4			
Via 1:	Ah	1	???				blt x5,x10,loop			
Via 2:	7h	1	fsw f1,0(x13)				fadd.s f1,f2,f1			
Via 3:	8h	1	addi x12,x12,4				addi x11,x11,4			
			7	6	5	4	3	2	1	0

Hits: 15
Misses: 18
Acessos: 33

Estrutura da cache

Após a 3ª iteração do loop:

Mapeamento direto

Via única:		TAG	V	DATA								
	0:	2h	1	addi x12,x12,4				addi x11,x11,4				Hits: 23 Misses: 23 Acessos: 46
	1:	2h	1	addi x5,x5,1				addi x13,x13,4				
	2:	2h	1	???				blt x5,x10,loop				
	3:	1h	1	fsw f1,0(x13)				fadd.s f1,f2,f1				
				7	6	5	4	3	2	1	0	

Completamente associativa

	TAG	V	DATA							
Via 0:	Ah	1	???				blt x5,x10,loop			
Via 1:	7h	1	fsw f1,0(x13)				fadd.s f1,f2,f1			
Via 2:	8h	1	addi x12,x12,4				addi x11,x11,4			
Via 3:	9h	1	addi x5,x5,1				addi x13,x13,4			
			7	6	5	4	3	2	1	0

Hits: 21
Misses: 25
Acessos: 46



Alinhamento das palavras em memória

Impacto no acesso à cache

RISC-V Instruction Set Architecture (ISA)

Dimensão das palavras na memória

...	
000Fh	Byte 15
000Eh	Byte 14
000Dh	Byte 13
000Ch	Byte 12
000Bh	Byte 11
000Ah	Byte 10
0009h	Byte 9
0008h	Byte 8
0007h	Byte 7
0006h	Byte 6
0005h	Byte 5
0004h	Byte 4
0003h	Byte 3
0002h	Byte 2
0001h	Byte 1
0000h	Byte 0

...	
000Fh	Half-Word
000Eh	7
000Dh	Half-Word
000Ch	6
000Bh	Half-Word
000Ah	5
0009h	Half-Word
0008h	4
0007h	Half-Word
0006h	3
0005h	Half-Word
0004h	2
0003h	Half-Word
0002h	1
0001h	Half-Word
0000h	0

...	
000Fh	Word 3
000Eh	/
000Dh	SP FP 3
000Ch	
000Bh	Word 2
000Ah	/
0009h	SP FP 2
0008h	
0007h	Word 1
0006h	/
0005h	SP FP 1
0004h	
0003h	Word 0
0002h	/
0001h	SP FP 0
0000h	

...	
000Fh	
000Eh	
000Dh	Double
000Ch	Word 1
000Bh	/
000Ah	DP FP 1
0009h	
0008h	
0007h	
0006h	
0005h	Double
0004h	Word 0
0003h	/
0002h	DP FP 0
0001h	
0000h	

RISC-V Instruction Set Architecture (ISA)

Alinhamento das palavras na memória

...	
000Fh	Byte 15
000Eh	Byte 14
000Dh	Byte 13
000Ch	Byte 12
000Bh	Byte 11
000Ah	Byte 10
0009h	Byte 9
0008h	Byte 8
0007h	Byte 7
0006h	Byte 6
0005h	Byte 5
0004h	Byte 4
0003h	Byte 3
0002h	Byte 2
0001h	Byte 1
0000h	Byte 0

O primeiro byte de cada Half-word é par, i.e., o bit menos significativo do endereço é zero.

000Bh	Half-Word
000Ah	5
0009h	Half-Word
0008h	4
0007h	Half-Word
0006h	3
0005h	Half-Word
0004h	2
0003h	Half-Word
0002h	1
0001h	Half-Word
0000h	0

O primeiro byte de cada Word é múltiplo de 4, i.e., os 2 bits menos significativos do endereço são zero.

000Bh	Word 2
000Ah	/
0009h	SP FP 2
0008h	
0007h	Word 1
0006h	/
0005h	SP FP 1
0004h	
0003h	Word 0
0002h	/
0001h	SP FP 0
0000h	

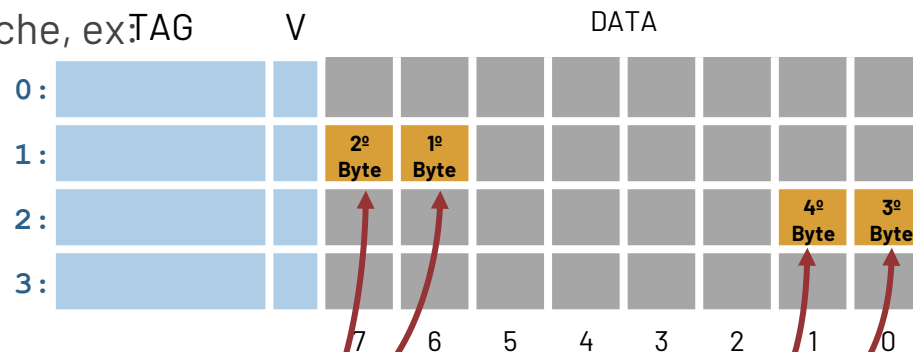
O primeiro byte de cada Double Word é múltiplo de 8, i.e., os 3 bits menos significativos do endereço são zero.

000Bh	/
000Ah	DP FP 1
0009h	
0008h	
0007h	Double Word 0
0006h	/
0005h	DP FP 0
0004h	
0003h	
0002h	
0001h	
0000h	

Acesso a variáveis não alinhadas em memória

- Considere um acesso não alinhado à cache, ex:TAG

lw x5,0 (X10)
Endereço base = 12Eh
↓ 4 bytes



Endereço do 1º Byte: 12Eh → TAG = 9, INDEX = 1, OFFSET = 6
Endereço do 2º Byte: 12Fh → TAG = 9, INDEX = 1, OFFSET = 7
Endereço do 3º Byte: 130h → TAG = 9, INDEX = 2, OFFSET = 0
Endereço do 4º Byte: 131h → TAG = 9, INDEX = 2, OFFSET = 1

- Se o endereço não estiver alinhado em memória é possível que a palavra esteja partida em duas linhas da cache!!!!
 - Resulta em dois acessos (cada um pode dar HIT ou MISS) em vez de um único acesso



Otimização de código para a cache

Se eu não quisesse desenhar processadores... qual a relevância das caches?

Exemplo de programa

Mesmo programa,
diferentes estruturas
de dados

Programa A

```
struct DATA{
    int a,b,c,d,e,f,g,h;
};

#define DATA_LEN
64*1024*1024
...

Struct DATA
myData[DATA_LEN];

...

for (i=0;i<DATA_LEN;i++) {
    myData[i].a = myData[i].b;
}
```

Programa B

```
struct DATA{
    int a, b;
};

#define DATA_LEN
64*1024*1024
...

Struct DATA
myData[DATA_LEN];

...

for (i=0;i<DATA_LEN;i++) {
    myData[i].a = myData[i].b;
}
```

Exemplo de programa

Mesmo programa,
diferentes estruturas
de dados

Programa A

```
struct DATA{
    int a,b,c,d,e,f,g,h;
};

#define DATA_LEN
64*1024*1024
...

Struct DATA
myData[DATA_LEN];

...

for (i=0;i<DATA_LEN;i++) {
    myData[i].a = myData[i].b;
}
```

176 ms

Programa B

```
struct DATA{
    int a, b;
};

#define DATA_LEN
64*1024*1024
...

Struct DATA
myData[DATA_LEN];

...

for (i=0;i<DATA_LEN;i++) {
    myData[i].a = myData[i].b;
}
```

142 ms

Intel i7-6700K
@4GHz

Programa A

Programa B

O bloco do processador é de 64 bytes, o que corresponde a 16 inteiros de 4B

Linha de cache no programa A: carregamento de 2 estruturas



Linha de cache no programa B: carregamento de 8 estruturas



O MISS RATE no acesso aos dados é:

caso A = $1/4 = 25\%$ caso B = $1/16 = 6,25\%$

Nota: o impacto na performance não é tão grande como seria de esperar devido à existência de um mecanismo de prefetch, que tenta adivinhar a sequência de acessos e faz pre-load das próximas linhas, e porque o Intel i7-6700K é um processador com execução fora de ordem, o que permite parcialmente esconder a latência dos acessos.

Intel i7-6700K
@4GHz

176 ms

142 ms

Programa C

```
#define ROW_SIZE 10*1024
#define COL_SIZE 10*1024

long long myMatrix[ROW_SIZE*COL_SIZE];

...
for (i=0 ; i < ROW_SIZE ; i++) {
    /* for each column */
    for (j=0 ; j < COL_SIZE ; j++) {
        /* for each row */
        myMatrix[i*COL_SIZE + j]++;
    }
}
```

Programa D

```
#define ROW_SIZE 10*1024
#define COL_SIZE 10*1024

long long myMatrix[ROW_SIZE*COL_SIZE];

...
for (j=0 ; j < COL_SIZE ; j++) {
    /* for each row */
    for (i=0 ; i < ROW_SIZE ; i++) {
        /* for each column */
        myMatrix[i*COL_SIZE + j]++;
    }
}
```

A matriz está organizada em memória guardando dois elementos vizinhos de uma dada linha em endereços adjacentes.

O elemento 11 é armazenado num endereço adjacente ao elemento 10

O elemento 28 é armazenado num endereço adjacente ao elemento 27

O elemento 48 é armazenado num endereço adjacente ao elemento 47

0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23
24	25	26	27	28	29	30	31
32	33	34	35	36	37	38	39
40	41	42	43	44	45	46	47
48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63

Pro

#d
#d

lo

..

fo

}

A matriz está organizada em memória guardando dois elementos vizinhos de uma dada linha em endereços adjacentes.

O programa é mais eficiente se efectuar as operações por:

[A] Colunas (programa C)

[B] Linhas (programa D)

[C] É indiferente

[D] Isto é tudo muito estranho

A	B							
	0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15	
16	17	18	19	20	21	22	23	
24	25	26	27	28	29	30	31	
32	33	34	35	36	37	38	39	
40	41	42	43	44	45	46	47	
48	49	50	51	52	53	54	55	
56	57	58	59	60	61	62	63	

];

A matriz está organizada em memória guardando dois elementos vizinhos de uma dada linha em endereços adjacentes.

O programa é mais eficiente se efectuar as operações por:

[A] Colunas (programa C)

Por cada acesso à memória são carregadas para a cache 7 operandos

Miss Rate = 12.5%

A

0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23
24	25	26	27	28	29	30	31
32	33	34	35	36	37	38	39
40	41	42	43	44	45	46	47
48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63

A matriz está organizada em memória guardando dois elementos vizinhos de uma dada linha em endereços adjacentes.

O programa é mais eficiente se efectuar as operações por:

[B] Linhas (programa D)

Por cada acesso à memória são carregadas para a cache 7 operandos

No entanto, como a matriz é de grandes dimensões, estes dados são retirados da cache antes de serem utilizados

Miss Rate = 100%

B							
0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23
24	25	26	27	28	29	30	31
32	33	34	35	36	37	38	39
40	41	42	43	44	45	46	47
48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63

Programa C

```
#define ROW_SIZE 10*1024
#define COL_SIZE 10*1024

long long myMatrix[ROW_SIZE*COL_SIZE];

...
for (i=0 ; i < ROW_SIZE ; i++) {
    /* for each column */
    for (j=0 ; j < COL_SIZE ; j++) {
        /* for each row */
        myMatrix[i*COL_SIZE + j]++;
    }
}
```

172 ms

Programa D

```
#define ROW_SIZE 10*1024
#define COL_SIZE 10*1024

long long myMatrix[ROW_SIZE*COL_SIZE];

...
for (j=0 ; j < COL_SIZE ; j++) {
    /* for each row */
    for (i=0 ; i < ROW_SIZE ; i++) {
        /* for each column */
        myMatrix[i*COL_SIZE + j]++;
    }
}
```

967 ms

Programa C	Programa D
<pre> #define ROW_SIZE 10*1024 #define COL_SIZE 10*1024 long long myMatrix[ROW_SIZE*COL_SIZE]; ... for (i=0 ; i < ROW_SIZE ; i++) { /* for each col */ for (j=0 ; j < COL_SIZE ; j++) { /* for each element */ myMatrix[i*COL_SIZE + j]++; } } </pre>	<pre> #define ROW_SIZE 10*1024 #define COL_SIZE 10*1024 long long myMatrix[ROW_SIZE*COL_SIZE]; ... for (j=0 ; j < COL_SIZE ; j++) { /* for each row */ for (i=0 ; i < ROW_SIZE ; i++) { myMatrix[i*COL_SIZE + j]++; } } </pre>

CONCLUSÃO:

É fundamental que o programador do software tenha consciência dos mecanismos de hardware do processador para que o programa desenvolvido tenha o melhor desempenho possível.

172 ms

967 ms