

MEMÒRIA CACHE

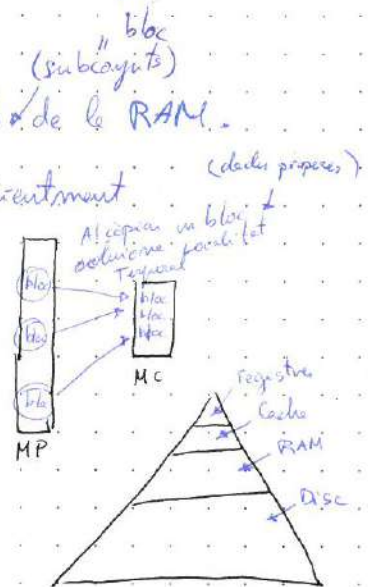
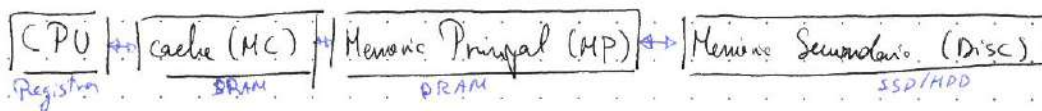
Introducció

SRAM = "Static". Es fa servir 6 Transistors (en DRAM 1) i així fa que sigui més cara.

Es fa servir com memòria cache. Més ràpida que DRAM.

Localitat Temporal = Es refereix que si una dada ha estat accedida una vegada, és més probable que es torni a accedir en un futur pròxim.

Memòria Cache: Situada entre CPU i RAM. Emmagatzema còpies temporals de la RAM. La seva funció és emmagatzemar dades que es faran servir freqüentment.



Terminologia

• **Fallada:** Primer la CPU llegirà la cache. Si no està en diu "fallada" (miss).

Després es copiarà de MP a MC el bloc de la dada.

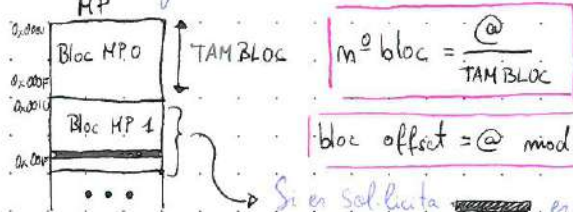
• **Reemplament:** Si el bloc que hem de copiar està ocupat amb un bloc antic, es reemplacen.

• **Emert:** La dada si que està en la MC.

Disseny bàsic d'una cache

Organització de la memòria en blocs

Cache guarda blocs, no dades "sueltas". Bloc té mida TAMBLOC bytes que són potències de 2.



$$n^o \text{ bloc} = \frac{\text{@}}{\text{TAMBLOC}}$$

Com que TAMBLOC és pot. de 2, la divisió correspon a desplaçant "eliminem el $\log_2(\text{TAMBLOC})$ bits de meigs per del n^o bloc".

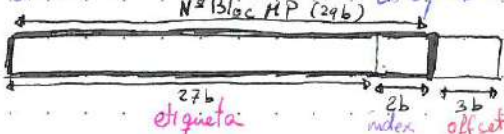
$$\text{bloc offset} = @ \bmod \text{TAMBLOC}$$

Això correspon a agafar els $\log_2(\text{TAMBLOC})$ bits meigs significatius de la @.

Si es sol·licita es copia tot.

Per guanyar velocitat (tot que no sigui eficient) sobre la velocitat que si $@ \% T = X$

es ve sobre el bit X. Exemple $0 \times 10010000 \rightarrow X=0$ i $0 \times 10010010 \rightarrow X=0$ en cas que $T=16$



etiqueta = Identificar forma única bloc de mem.

index = Det. quina línia de cache es farà servir per aquests blocs.

offset = Identificar part específica dins el bloc.

Procediment

1) S'extreu index de la direcció de memòria.

2) Es fa servir index per acudir línia en cache.

3) Comparar etiq. de la direcció mem amb l'etiq. guardada en aquella línia de la cache.

4) if (etiq ==) Hit ; else Miss, carreguem nou bloc sobreescribant.

Exemple: $0x00000028 = (0000.00101000)_2$ en Bloc 2^3 $\boxed{\#Bloc MP=5}$ i 2^2 línies $\boxed{\#Línia=1}$

00101000 eliminem els 3 últims i el m^e que queda $(0000.0101)_2 = (5)_{10}$ i aquest és el m^e bloc.

00101000 al fer mòdul 4 agafem 2 últims (abans fem div) i aquesta és la línia en Cache.

0000.00101000 aquesta és l'etiqueta.

Exemple memòria cache

V	Etiqueta	Word 1	Word 2
1	0x0	0x00000000	0x00000000

#La cache en aquest exemple és de 8B per això té Word 1, Word 0.

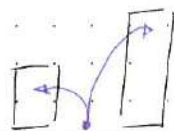
* Si la grandària del bloc de cache és massa gran hi haurà molts miss (Donat que no hi haurà molts blocs al ser SRAM car) i cada miss s'haurà de fer copy gran i això farà process lent.

#64Kb és el m^e optm (en 2012).

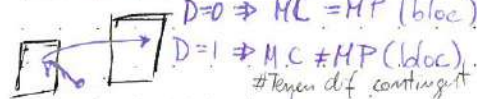
Gestió de les escriptures

Esriptura immediat (Write-through) [Enort]

S'escriu simultàniament a la memòria cache i a la principal.



Esriptura retardada (Write-back o Copy-back) [Enort]



S'escriu només a la cache. S'escriu a la MP quan la línia de la cache serà sobreescrita.

Esriptura amb assignació (Write allocate) [Miss]

Si es produeix miss, copiem el bloc a la MC (igual que lectura).

Esriptura sense assignació (No-Write allocate) [Miss]

Si es produeix miss, modifiquem contingut en MP, com si no hi hagués cache.

• Esriptura immediat amb assignació (MARS)

[Hit]: Escriu a MP; MC a la vegada.

[Miss]: Còpia bloc MP \rightarrow MC i després escriu MP; MC a la vegada.

• Esriptura immediat sense assignació

[Hit]: Escriu a MP; MC a la vegada.

#Recorda que SEMPRE MP \rightarrow MC en cas de miss en lectura.

[Miss]: Només escriu en bloc de MP (hi haurà dif. contingut MP \rightarrow Nou, MC \rightarrow Vell).

• Esriptura retardada amb assignació

[Hit]: Només s'escriu en MC i Bit Dirty = 1. (Hi haurà dif. cont. MP \rightarrow Vell, MC \rightarrow Nou).

[Miss]: Store \rightarrow Còpia el bloc MP \rightarrow MC i modifiquem MC: D = 1.

Load \rightarrow Còpia MC \rightarrow MP \leftarrow i després el bloc de MP \rightarrow MC. Com que és Load D = 0.

Medida de rendimento

$$t_{\text{acc}} = t_h + t_p$$

t_{acc} = Tempo de servir d'una referència a memòria

t_h = Temps per determinar si es tracta o fallada

t_p = Temps penalització per resoldre referència

Buffer d'Escriptura: Quan tinguem escriptura immediata. Suponem que no caldrà esperar a MP.

T_p	Immediata sense assig	Ref. amb assig
Lectura - Éxit	0	0
Lectura - Fallada	$t_{\text{block}} + t_h$	bloc mod: $2t_{\text{block}} + t_h$ bloc no mod: $t_{\text{block}} + t_h$
Escriptura - Éxit	0	0
Escriptura - Fallada	0	bloc mod: $2t_{\text{block}} + t_h$ bloc no mod: $t_{\text{block}} + t_h$

2 pg. h. de penat. d'unitat a MP.

$$t_p = (1 - p_e)(t_{\text{block}} + t_h)$$

proporció escriptures

$$t_p = p_m(2t_{\text{block}} + t_h) + (1 - p_m)(t_{\text{block}} + t_h)$$

proporció modificats

no time p.e. per a MP h. de penat. entre escriptura i lectura

$$t_{\text{ma}} = t_h + m * t_p$$

taxa miss

$$M_{\text{ins}} = \frac{M_{\text{r}}}{I}$$

núms referències per instrucció

$$t_{\text{ma}} = \frac{t_{\text{ma}} : \text{fèlta} + (M_{\text{ins}} - 1) * t_{\text{ma}} : \text{dada}}{M_{\text{ins}}}$$

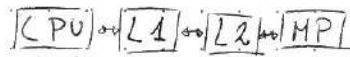
en aquest cas és mitjà

temp. acte

$$t_{\text{ee}} = M_{\text{ins}} * \text{CPI} * t_c = (M_{\text{ins}} * \text{CPI}_{\text{ideal}} + M_{\text{falledes}} * t_{\text{pm}}) * t_c$$

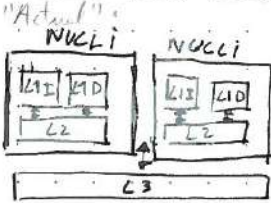
Multimèdia

V1:

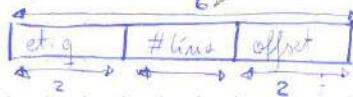


L1: Blocs petits, molt ràpids. To. gran

L2: Blocs més grans, menys ràpids (Tall. fèr. per no anar a MP). Contens.



6.1. $HP = 64B$ $MC = 16B$ $1B_{loc} = 4B$ | $MC = 2 \text{ línies} \cdot 2 \text{ byte/blos}$



@ $0x39 = 00111001$

OFFSET = $01_2 = 1_{10}$
#Línia = $10_2 = 2_{10}$
#Etq = $1110_2 = 14_{10}$

#Bloc HP	binari	#línia MC	hit/miss
0	0000	0	m
1	0001	1	m
0	0000	0	H
3	0011	3	m
12	1100	0'	m
10	1010	2	m
3	0011	3	H
7	0111	3'	m
12	1100	0'	H
15	1111	3''	m

#pg. està buit
#pg. està buit
#pg. està buit
el cont no correspon
#pg. està buit
#pg. cont. no correspon
#pg. cont. no correspon

6.2. 16b @. Mida total = $256B$ Mida Bloc = $16B$ Escripció immediata sense assign.

a) Omple la Taula

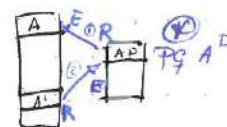
tipus	etq (Hexa)	índex MC	H/M	#B lleg HP	#B esc HP	Lec. d'adreça en MC	Esc. d'adreça en MC
R	4534	0x3	M	16	0	NO	SI
R	4568	0x6	M	16	0	NO	SI
W	13A4	0xA	M	0	16	NO	NO
W	13A8	0xA	M	0	2	NO	NO
R	3560	0x6'	M	16	0	SI	SI
W	453C	0x3	M	0	2	NO	SI
W	60A0	0xA'	F	0	2	NO	NO
R	453C	0x3	E	0	0	SI	NO
W	3900	0x0	F	0	2	NO	NO
R	A238	0x3'	F	16	0	SI	SI

índex i índex però dif. etq.

No form. fàcil
Pg. work
3on 2B al 16b
#línia
offset
Pg. el bit és real amb el que hi ha a MC
#Tercer bit
immediat
memòria mod HP
no pota el mod a MC
Carrega bloc a MC
Demo resultat desde MC

b) Omple amb | Esc. Ret. amb assign.

R	4534	0x3	F	16	0	SI	SI
R	4568	0x6	F	16	0	SI	SI
W	13A4	0xA	F	16	0	NO	SI
W	13A8	0xA	E	0	0	NO	SI
R	3560	0x6'	F	16	0	SI	SI
W	453C	0x3	E	0	0	NO	SI
W	60A0	0xA'	F	16	0	SI	SI
R	453C	0x3	E	0	0	SI	NO
W	3900	0x0	F	16	0	NO	SI
R	A238	0x3'	F	16	16	SI	SI



amb amic (ata mod en MC)
Retardable (mod de mod)
Retardable
Assign + Register
Amb amic
línia()

F. xet que el read és de 3^a així que hem de portar el vell a HP abans de portar HP a MC i després netejar el MC.

6.7) CPI = 1.5 cycles/ins. // $t_c = 10 \text{ ns}$ // $nr = 1.6$ // $\begin{matrix} \text{dada} \\ \text{Ins.} \end{matrix}$ // $\begin{matrix} \text{exc.} \\ \text{ret.} \\ \text{amb. ang.} \end{matrix}$ // $T_{amb.}$

a) Quin serà t_{amb} (Teix amb mem. mitja) en cycles?

$$t_{ma} = t_h + nr * t_p$$

$$t_{ma} = \frac{t_{ma:f} + (nr_{ins} - 1) * t_{ma:dada}}{nr_{ins}}$$

$$t_{ma:f} = t_{h:f} + m_f * t_{p_f} \quad t_{p-mod} \quad t_{p-nomod}$$

$$t_{ma:d} = t_{h:d} + m_d * (p_m (2 * t_{p_{mod}} + t_h) + (1 - p_m) (t_{p_{nomod}} + t_h))$$

$$t_{ma:f} = 1 + 0.04 * 10 = 1 + 0.4 = 1.4$$

9 no està en sol?

$$t_{ma:d} = 1 + 0.1 * (0.2 (2 * 20 + 1) + (1 - 0.2) (15 + 1)) = 1 + 0.1 (0.2 * 41 + 0.8 * 16) = 1 + 0.1 (8.2 + 12.8) = 1 + 0.1 (21) = 1 + 2.1 = 3.1$$

$$t_{ma} = \frac{1.4 + (1.6 - 1) * 3.1}{1.6} = \frac{1.4 + 0.6 * 3.1}{1.6} = \frac{1.4 + 1.86}{1.6} = \frac{3.26}{1.6} = 2.03 \text{ cycles}$$

b) I t_{exec} en ns?

$$t_{exec} = (nr_{ins} * CPI_{ideal} + m_f * t_{p_f} + m_d * t_{p_d}) * t_c$$

$$t_{exec:ins} = CPI * t_c$$

$$CPI = CPI_{ideal} + penal_f + penal_d$$

$$penal_f = m_f * t_{p_f}$$

$$penal_d = (nr - 1) * m_d (p_m * t_{p_{mod}} + (1 - p_m) t_{p_{nomod}})$$

$$penal_f = 0.04 * 10 = 0.4$$

$$penal_d = (1.6 - 1) * 0.1 (0.2 * 20 + (0.8 * 15)) = 0.6 * 0.1 * (16) = 0.06 * 16 = 0.96$$

$$CPI = 1.5 + 0.46 + 0.4 = 2.36$$

$$t_{exec} = 2.36 * 10 * 10^{-9} = 2.36 * 10^{-8} \text{ s}$$

EC Examen de Problemes

Exercici 1 (Examen Final juny 2011)

Considera un computador amb un processador que té amplada de dades i d'adreces de 64 bits, i una memòria cache de dades amb les següents característiques:

- 512 blocs, amb 2 paraules per bloc (paraules de 64 bits) $1024 \text{ words} = 8192$
- correspondència directa
- escriptura immediata sense assignació

- Quina és la capacitat en bytes per a dades de la memòria cache? 8192 Bytes
- Indica el rang de bits de l'adreça que especifiquen l'índex a la memòria cache $[4, 12]$
- Indica el rang de bits de l'adreça que especifiquen l'etiqueta $[13, 63]$
- Quants bits d'emmagatzematge per a etiquetes i bits de control fan falta en total per cada entrada de la memòria cache? $51b$
- Considera el següent programa en alt nivell, que s'executa en aquest computador:

```
int V[6];          /* un int ocupa 64 bits */
main() {
    int i, tmp;      /* variables ubicades en registres */

    tmp = V[0];
    for (i=1; i<6; i++) V[i-1] = V[i];
    V[5] = tmp;
}
```

@	V/W	H/M
0	r	m
4	r	m
8	w	h
12	r	m
16	w	h
20	r	m
24	w	h
28	r	m
32	w	h

Tenint en compte que el vector V està emmagatzemat a partir de l'adreça 0, indica la seqüència d'adreces (en hexadecimal) dels accessos a memòria de dades que genera l'execució del programa, especificant per cada una: si és lectura o escriptura (R/W) i si produeix un encert o fallada (hit/miss) a la cache.

Exercici 2 (Examen Final gener 2013)

Considera el següent programa:

```
int M[4][2];      /* adreça base d'M = 0 */
int V[2];
main() {
    int i, j;      /* emmagatzemades en registres */
    for (i=0; i<4; i++)
        for (j=0; j<2; j++)
            M[i][j] = M[i][j]+V[j];
}
```

que s'executa en un computador MIPS que disposa d'una memòria cache de dades, inicialment buida, de correspondència directa i política d'escriptura retardada amb assignació, que conté 4 blocs i on els blocs són de 8 bytes.

Emplena la següent taula, que mostra la seqüència de les 12 primeres referències a memòria (E: escriptura/ L: lectura) corresponent al programa.

0
1
2
3

- (1) Entre de MP-DML a mem MC.
 (2) Com que en VL17 capta VL0; VL1 i la resta D=1
 (3) Només mod MC.

87 fixat que un bloc són 8B.
 Quan fa L1MZO[CO] al bit 4B, capta
 també MZO[CO] a 6MC.

No estue pq és línies 0 i no 1 (Hem de copiar VL1; VL2) tot que no existeixi no?
 No estue pq 8B de l'2

	element accedit	línia de MC	hit/miss	bytes llegits d'MP	bytes escrits a MP
L	M[0][0]	0	miss	8	0
L	V[0]	32%4=0	miss	8	0
E	M[0][0]	0	miss	8	0
L	M[0][1]	0	hit	0	0
L	V[1]	2	miss	8	0
E	M[0][1]	2	miss	8	0
L	M[1][0]	1	miss	8	0
L	V[0]	0	miss	8	8
E	M[1][0]	1	hit	0	0
L	M[1][1]	1	hit	0	0
L	V[1]	0	hit	0	0
E	M[1][1]	1	hit	0	0

Table 1: TLB

Exercici 3 (Examen Final juny 2012)

Considera el següent programa:

```

int M[F][C];          /* adreça base d'M = 0 */

main() {
  int i,j;              /* emmagatzemades en registres */

  for (i=0; i<F-4; i++)
    for (j=0, j<C; j++)
      M[i][j] = M[i+4][j];
}
  
```

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

que s'executa en un computador MIPS que disposa d'una memòria cache de dades, inicialment buida, de correspondència directa, que conté 16 blocs i on els blocs són de 16 bytes.

Omple una taula especificant el nombre de referències, el nombre de fallades i el nombre de bytes transferits a/des de MP per aquests dos casos:

- MC amb política d'escriptura immediata sense assignació; matriu M amb F=16 i C=16
- MC amb política d'escriptura retardada amb assignació; matriu M amb F=8 i C=8

Esc. imm. sense assign. C=8 F=16 " C=16

	W/R	h/m	#r	#w
* H[4][0]	R	m	16	0
H[0][0]	W	m	0	16
H[4][1]	R	h	0	0
H[0][1]	W	m	0	4
H[4][2]	R	h	0	0
H[0][2]	W	m	0	4
H[4][3]	R	h	0	0
H[0][3]	W	m	0	4
* H[4][4]	R	m	16	0
H[0][4]	W	m	0	16
H[4][5]	R	h	0	0
H[0][5]	W	m	0	16
H[4][6]	R	h	0	0
H[0][6]	W	m	0	16

$$16B * (12 * 16) = 3072$$

$$N^{\circ} \text{ refs} = 192 = 12 * 16$$

$$N^{\circ} \text{ ref} = 2 * 192 = 384 = N^{\circ} \text{ ref}$$

$$Q < 1644 \Rightarrow 12 \text{ línies}$$

$$Q < 16 \text{ pero 4 repes per file}$$

$$12 * 4 * (5) = 240 \text{ misses}$$

$$12 * 4 * (16 + 4 + 4 + 4) = 1536B$$

Esc. retardada amb assign. F=8 " C=8

	W/R	h/m	#r	#w
R M[4][0]	m	16	0	
W M[0][0]	m	16	0	
R M[4][1]	h	0	0	
W M[0][1]	h	0	0	
R M[4][2]	h	0	0	
W M[0][2]	h	0	0	
R M[4][3]	h	0	0	
W M[0][3]	h	0	0	
R M[4][4]	m	16	0	
W M[0][4]	m	16	0	
R M[4][5]	h	0	0	
W M[0][5]	h	0	0	
R M[4][6]	h	0	0	
W M[0][6]	h	0	0	

$$N^{\circ} \text{ ref} = (8 * 8) * 2 = 64 * 2 = 128 = N^{\circ}$$

Mi hanie 16 miss pq per cada bloc hi ha miss i no hi ha writeback

$$16 * 16 = 256B$$

Ens fixem que tot i que hi hagi Dirty a 1 en alguns blocs men hi ha sobreescritura.

EC Examen de Problemes

Exercici 1 (Examen Final 2011-2012 Q1)

Suposem que tenim un processador de 32 bits amb una memòria cache de dades de 8 KB associativa per conjunts de 2 vies, on cada bloc té 16 bytes, i que es segueix l'algorisme de reemplaçament LRU.

- Calcula el nombre de fallades de la cache en executar el següent programa, suposant que la cache té la política d'escriptura immediata sense assignació, i que la memòria cache és inicialment buida. L'adreça base del vector A és 0.

```
int A[1024]; 0 -> 4096  A0 = 0000.0000.0000.0
int B[1024]; 4096 -> 8192  B0 = 0000.0000.0000.1
int C[1024]; 8192 -> 12288  C0 = 0000.0000.0000.0
```

```
void main() {
    int i;
    for (i=0; i<1024; i++)
        A[i] = B[i] + C[i];
}
```

falladesA= 1024
falladesB= 256
falladesC= 256

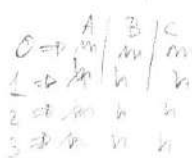
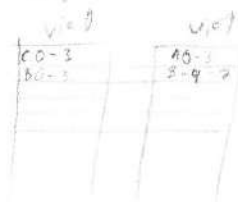
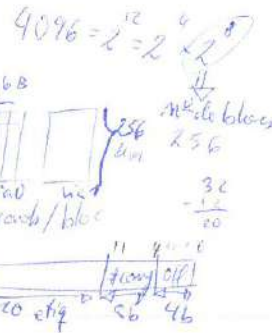
- Fes el mateix per al següent programa:

```
int A[512]; 0 -> 2048  A0 = 0000.0000.0000.1
int B[512]; 2048 -> 4096  B0 = 0000.0000.0000.0
int C[512]; 4096 -> 8192  C0 = 0000.0000.0000.1
```

```
void main() {
    int i;
    for (i=0; i<512; i++)
        A[i] = B[i] + C[i];
}
```

falladesA= 512
falladesB= 128
falladesC= 128

- Repeteix els dos apartats anteriors considerant ara que la cache té una política d'escriptura retardada amb assignació.

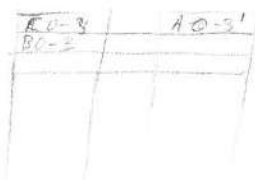


al ser sense assign. A no es guarda i B i C s'han de carregar cad un.

A cada iter. es van sobre escriure.

falladesA= 1024
falladesB= 1024
falladesC= 1024

falladesA= 128
falladesB= 128
falladesC= 128



Exercici 2 (Examen Final 2013-2014 Q1)

Suposem que tenim un processador de 32 bits amb una memòria cache de dades de 256 bytes, on cada bloc té 16 bytes. Suposem que executem els següents programes.

```
//programa A
int M[4][64];
```

```
void main() {
    int i, j; //en registres
    for (i=0; i<4; i++)
        for (j=0; j<64; j++)
            M[i][j] = 0;
}
```

```
//programa B
int M[4][64];
```

```
void main() {
    int i, j; //en registres
    for (j=0; j<64; j++)
        for (i=0; i<4; i++)
            M[i][j] = M[i][j] + 1;
}
```

$M[0][0] \rightarrow M[0][15]$ (m, h, h, h)
 $M[0][16] \rightarrow M[0][31]$ (m, h, h, h)

$M[0][0] = 0 \times 0$
 $M[1][0] = (1 \times 64 + 0) \times 4 = 256 \quad 0 \times 100$
 NO sobre escriu però cont. amb motes
 patró. Ve a memòria inferior.

$M[0][0] = M[0][0] + 1 \quad (0 \times 0)$
 $M[1][0] = M[1][0] + 1 \quad (0 \times 100)$
 $M[2][0] = M[2][0] + 1 \quad (0 \times 200) \quad (2 \times 64 + 0) \times 4 = 512$
 $M[3][0] = M[3][0] + 1 \quad (0 \times 300)$
 $M[0][1] = M[0][1] + 1$

Sobre escriu tota l'adreça

Calcula el nombre de fallades de la cache suposant que la memòria cache és inicialment buida. L'adreça base de la matriu M és 0.

- Suposant que la cache és de correspondència directa i té la política d'escriptura retardada amb assignació.

falladesA= 64 falladesB= 256 (4x64)

- Suposant que la cache és associativa per conjunts de 4 vics (algorisme de reemplaçament LRU), i que té la política d'escriptura immediata sense assignació.

falladesA=

falladesB=

256
mai hi ha
lectura
sempre miss

64 Del 0 a 3 (m, h, h, h) però després
sobre escriu
(64x4) = 256 però memòria ja té miss 256 (64)

0.0 a 0.3 (m, h, h, h)

	Vic 0	Vic 1	Vic 2	Vic 3
0.0 a 0.3	0.0 a 0.3	1.0 a 1.3	2.0 a 2.3	3.0 a 3.3
1.0 a 1.3				
2.0 a 2.3				
3.0 a 3.3				

4 blocs

1w = 4B

bloc 4 words
#index 3
offset 16



0000.0000	0000	#0
0001.0000	0000	#16
0010.0000	0000	#32
0011.0000	0000	#48