

## 3.8 Replacement policies

21.05.2019

a) Random

b) FIFO

c) LRU - Least Recently Used

No. of "age" bits =  $\log_2$  (set associativity)

Example 1: P. că avem un set asociativ pe 4 căi cu adrese în zecimal și pp. că toate adresele se mapează la același index

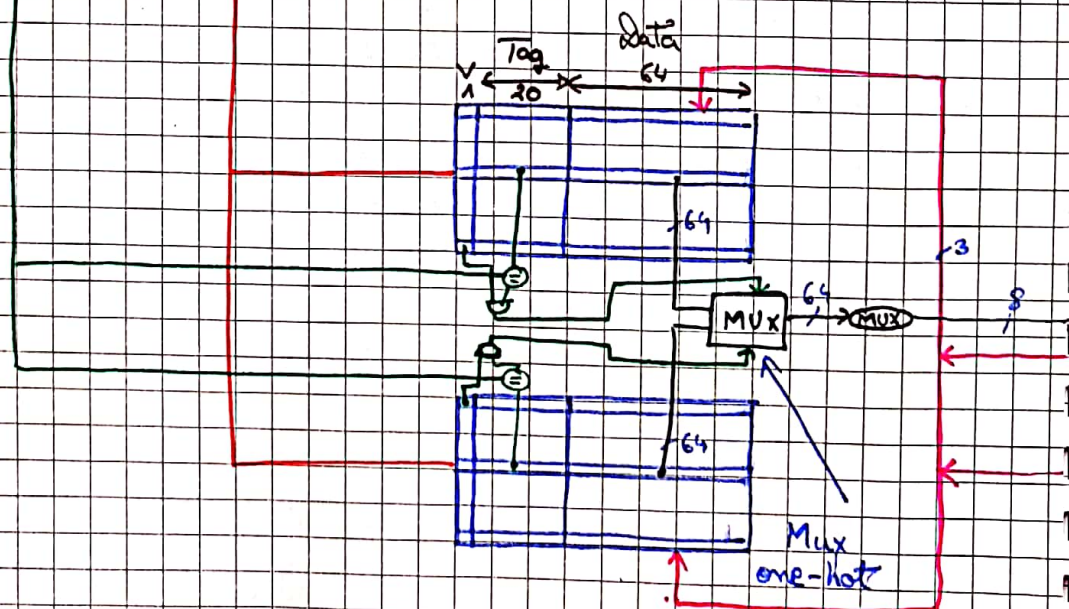
Tip  
Adresa de fapt  
e curent de la  
adresa în "mem" 14

		Date				Age			
		0	1	2	3	0*	1*	2*	3*
Miss	14	0	0	0	0	0	0	0	0
(cache gol)	14	0	0	0	0	1	0	0	0
M → 6	14	6	0	0	0	2	1	0	0
M → 12	14	12	0	0	0	3	2	1	0
M → 10	14	10	0	0	0	0	3	2	1
(Hit) H → 6	10	6	0	0	0	1	3	0	2
H → 7	10	7	0	0	0	2	0	1	3
M → 14	10	14	0	0	0	3	1	2	0
H → 7	10	7	0	0	0	3	0	2	1
M → 3	3	3	0	0	0	0	1	3	2
H → 14	3	14	0	0	0	1	2	3	0

← lei cu vârsta mai mare rămân la fel, cei cu mai mică se incrementează

Example 2: Memoria Vax 11 / 480 - 2-way SA

Tag	Index	B. off.
<20>	<3>	<3>





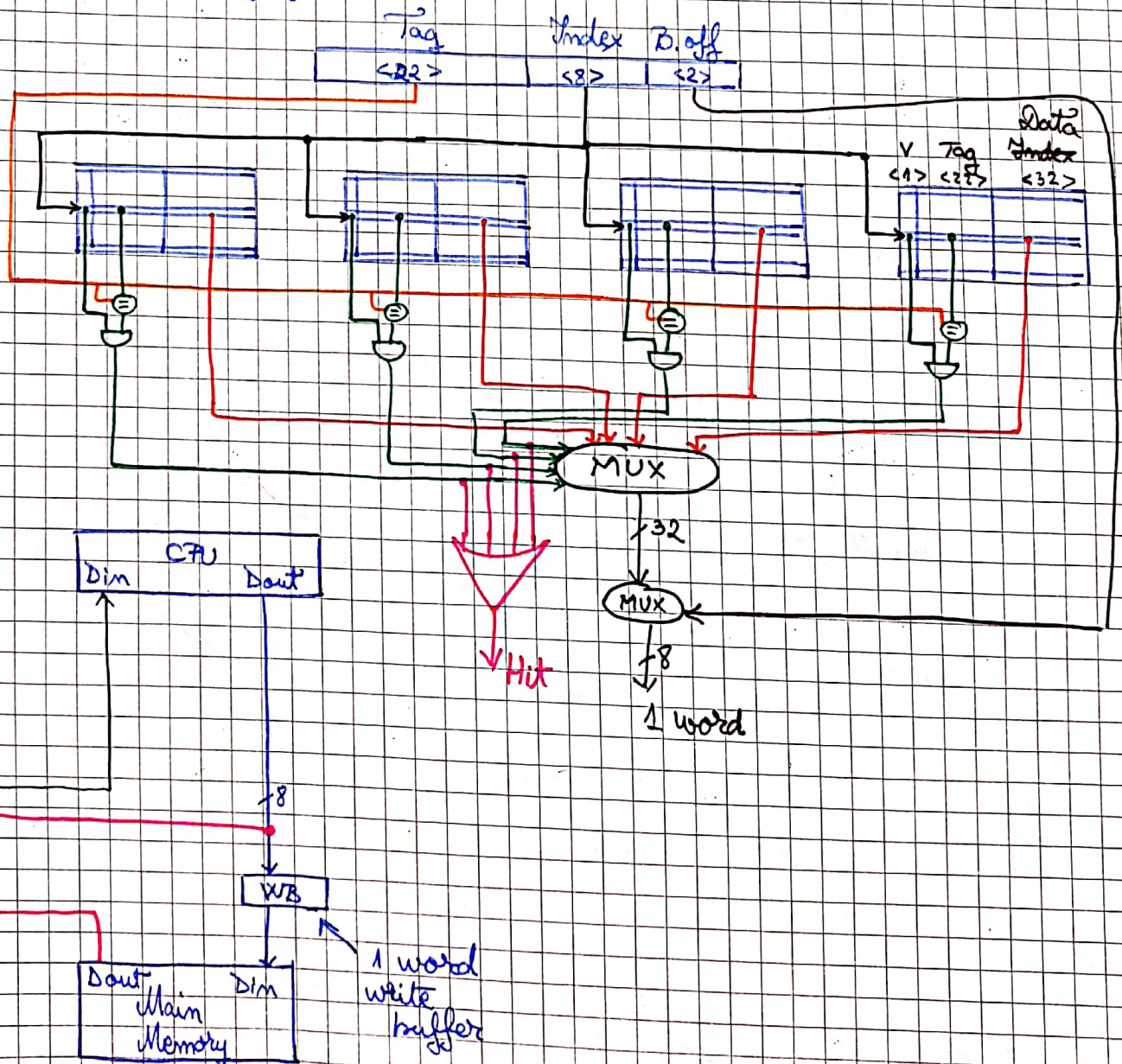
$2^3 \text{ words/block} = 2^3 \text{ B/block} = 2^3 \cdot 2^3 \text{ bits/block} = 64 \text{ bits/block}$   
 Cache Size (data) = p.a  $\times 2^{\text{index}}$   $\times$  block size =  $2 \times 2^3 \times 2^3 \text{ B} = 2^{13} \text{ B}$   
 $\uparrow$  set associativity  $= 8 \text{ KiB}$

Write Policy  $\rightarrow$  WriteThrough (Write Buffer)  $\rightarrow$  do Allocate in case of a Write miss.

Example 3:

- 4 way - S.A. ; 1 word = 1 B ; 1 block = 4 words =  $2^2$  words
- Cache size =  $2^{10}$  blocks = 1 KiBlocks
- MM size = 4 GiB

$$\frac{\text{No. of blocks}}{\text{p.a}} = \frac{2^{10}}{2} = 2^8 \text{ blocks/bank}$$





Example 4: Se dau 2 memorii cache diferite alor 2 sint. de calcul de ~~ant~~ altfel identice.

cache A: 2 way s.a ;  $I_{mr} = 2\%$  ;  $D_{mr} = 1,7\%$

cache B: 4 way s.a ;  $I_{mr} = 1,6\%$  ;  $D_{mr} = 1,4\%$

A și B sunt maximi de calcul load/stora, 20% din instrucțiuni fiind load/stora.

Absolute miss penalty = 200 ns

$CCT_A = 20$  ns

La B, pt. sa s.a e mai mare, mux e mai complicat, timpul e mai mare, deci  $CCT_B$  se mărește cu 10% (s-a degradat)

$CPI_{ideal} = 3$  clock cycle-wr

Care maxima e mai buna?

$$CPU_{time} = IC + \underbrace{\left( CPI_{ideal} + \text{Memory access per instr} \times \frac{\text{Miss Rate}}{\text{Miss Penalty}} \right)}_{\text{Misses per instr}} \times CCT$$

$$\text{Misses per instr. A} = 1 \times I_{mr} + 0,2 \times D_{mr} = 0,02 + 0,2 \times 0,017 = 0,0234$$

20%

$$\text{Misses per instr. B} = 0,016 + 0,2 \times 0,014 = 0,0188$$

$$\text{Miss Penalty A} = \left[ \frac{200 \text{ ns}}{20 \text{ ns}} \right] = 10 \text{ clock cycles}$$

$\uparrow$   
 $CCT_A$

$$\text{Miss Penalty B} = \left[ \frac{200 \text{ ns}}{22 \text{ ns}} \right] = 10 \text{ clock cycles (treb. sa fie mereu nr. intreg)}$$

$\uparrow$   
 $CCT_B = CCT_A + 20\% \cdot CCT_A$

$$CPU_{time A} = IC \times (3 + 0,0234 \times 10) \times 20 \text{ ns} = IC \times 3,234 \times 20 = IC \times 64,64 \text{ ns}$$

(timpul mediu pt. o instructiune, luand in considerare imperfecțiunea cache-ului)

$$CPU_{time B} = IC \times (3 + 0,0188 \times 10) \times 22 \text{ ns} = IC \times 70,136 \text{ ns}$$



Exemple 5: Admitem un cache de 64 kBytes:

$$CCT = 20 \text{ ns}$$

Media nr. de accese la instructiune: Memory access per. = 1,3  
instr.

$$CPI_{ideal} = 1,5$$

Pt. acest cache, daca se adopta set asociativ, det. o degradare a frecv. clock-ului cu 8,5%. In acest context, se analizeaza cache-ul in felul urm.: mapare directa vs. 2 way S.A.

$$\text{Statistic: Miss rate}_{DM} = 3,3\%$$

$$\text{Miss rate}_{2 \text{ way S.A.}} = 3\%$$

Care solutie e mai ok stiind ca in ambele cazuri, pt. Miss Penalty = 200 ns (expresie temporală).

In mod normal, Miss Penalty e exprimat in clockuri

$$AMAT = CCT + \text{Miss Rate} \times \frac{\text{Miss Penalty}}{\text{Time}}$$

$$CPU_{time} = \gamma C \times \left( CPI_{ideal} \times CCT + \text{Mem. acc per instr.} \times \text{Miss Rate} \times \frac{\text{Miss Penalty}}{\text{Time}} \right)$$

Miss Penalty clockes  $\times CCT$

$$AMAT_{DM} = 20 \text{ ns} + 0,033 \times 200 \text{ ns} = 20 \text{ ns} + 6,6 \text{ ns} = 26,6 \text{ ns}$$
$$= 20 + 33\% \cdot 20$$

Better

$$AMAT_{SA2} = 21,7 + 0,03 \times 200 \text{ ns} = 24,7 \text{ ns}$$

Better

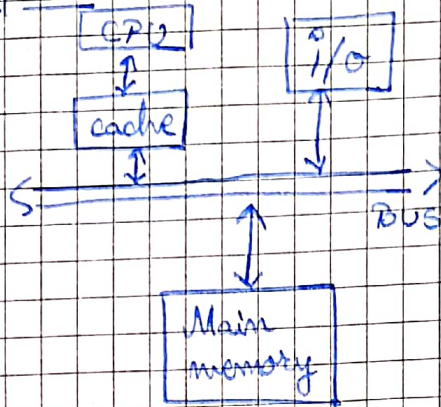
$$CPU_{time_{DM}} = \gamma C \times (1,5 \times 20 \text{ ns} + 1,3 \times 0,033 \times 200) \approx \gamma C \times 40,14 \text{ ns}$$

$$CPU_{time_{2SA}} = \gamma C \times (1,5 \times 21,7 \text{ ns} + 1,3 \times 0,03 \times 200) \approx \gamma C \times 40,35 \text{ ns}$$

$CPU_{time}$  ia in considerare mai mult ca AMAT,  
deci DM e mai bun.



### Example 6:



Cât la sută din BUS Bandwidth este „mâncat” de imperfecțiunea cacheului?

Se da un sist. de calcul cu:

$$\text{BUS Bandwidth} = 10^9 \text{ words/sec}$$

$$\text{Tp. ca: Miss Rate} = 10\%$$

$$1 \text{ block} = 4 \text{ words}$$

$$\text{CPU reference frequency} = 10^8 \text{ words/sec (reads \& writes)}$$

$$\text{BUS width} = 2 \text{ words}$$

P. ca în (t) mom., 35% din blocurile cacheului au fost modificate (blockuri dirty). Cache are o politică write-allocate, frecvența scrierilor e de 30%

Alternative de design, știind că Bandwidth e „mâncat”:

a) write bank (WB)

b) write time (WT)

a) Nb. de accese la Bus datorită imperfecțiunii cache =  $10^8 \times 0,1 \times (\% \text{ Reads} \times \text{Read Miss Penalty} + \% \text{ Writes} \times \text{Write Miss Penalty})$

$$\text{Read Miss Penalty}_{WB} = \left( \frac{\text{Block size}}{\text{Bus width}} \right) \times \text{Bus writes} \times 0,35 +$$

↳ nr. de accese pe bus

↳ blockuri dirty

$$+ \left( \frac{\text{Block size}}{\text{Bus width}} \right) \times \text{Bus reads}$$

Update

Allocate

$$= 2 \times 0,35 + 2 = 2,7 = \text{Write miss penalty}$$

$$\% \text{ Bus Used}_{WB} = \frac{10^8 \times 2,7}{10^9} = \frac{2,7}{10} = 2,7\%$$

b) Jemă: WT



$$\text{Nr. access} \dots = 10^8 \times 0,1 + (\% \text{ Reads} \times \text{Read Miss Penalty} + \% \text{ Writes} \times \text{Write Miss Penalty}) + 10^8 \times 0,9 \times 0,3 \times \text{Write Hit Penalty}$$

$\underbrace{\hspace{10em}}_{\text{Writes}} \quad \underbrace{\hspace{10em}}_{= 1}$

$$\begin{aligned} \text{Read Miss Penalty} &= \text{Write Miss Penalty} = \\ &= \underbrace{2 \text{ BUS Reads}}_{\text{Allocate}} + 0,35 \cdot \underbrace{2 \text{ BUS Writes}}_{\text{Update for Main Memory}} \end{aligned}$$

CVRS 13

28.05.2019

$$\begin{aligned} \text{b) Nr. access} &= \underbrace{10^8 \text{ words/second}}_{\text{Frequency of memory references}} \times \underbrace{0,1}_{\text{Miss Rate}} \times (0,3 \text{ Write Miss Penalty} + \\ &+ 0,7 \text{ Read Miss Penalty}) + 10^8 \times 0,9 \times 0,3 \text{ Write Hit Penalty} \end{aligned}$$

$$\text{Read Miss Penalty} = 2 \text{ BUS Reads (nur main update!)} = 2$$

$$\begin{aligned} \text{Write Miss Penalty} &= 2 \text{ BUS Reads} + 1 \text{ BUS Write} = \\ &= 3 \text{ BUS Accesses} \end{aligned}$$

$$\text{Write Hit Penalty} = 1 \text{ BUS Write} = 1 \text{ BUS Access}$$

$$\begin{aligned} \text{Adapt: Nr. access} &: 10^8 \cdot 0,1 (0,3 \cdot 3 + 0,7 \cdot 2) + 10^8 \cdot 0,9 \cdot 0,3 \cdot 1 = \\ &= 10^7 \cdot 2,3 + 2,7 \cdot 10^7 = 5 \cdot 10^7 \end{aligned}$$

$$\% \text{ Bus Used}_{WT} = \frac{10^7 \cdot 5}{10^9} = \frac{5}{100} = 5 \%$$

### 3.9. Reducing miss rate with multi-level caches

