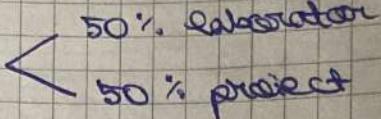


21.02.2023

ORGANIZAREA CALCULATORULOR

Curs 1

Exam 67%

Proiect 33% 

La examen : probleme (open book) 
(max 5 pb)
2p / problema

La laborator : 2 teste (pb. de antrenament)

La proiect :

LD \rightarrow AC \rightarrow OC \rightarrow FIC

Computer science

Computer engineering

Software Engineering

Information Technology

- performanta alg

- memorii

CPU time = IC * CPI * Clock cycle time

Power Wall

Networks - On - Chip

Federated Learning

Joi de la 10 în 18:30 Consultări

Capitel 1:

Booth's algorithm

$$\begin{array}{r} 125 \\ \boxed{7} \\ 0125 \downarrow \\ \text{gummi} \end{array}$$

$$\begin{array}{r} 0.12 \boxed{7} \\ \downarrow \\ 0.120 \end{array} \quad \xrightarrow{\text{gummi}}$$

$$Y = X$$

$$X = x_{m-1}x_{m-2}\dots x_1 \dots x_0 \rightarrow x_i \in B = \{0, 1\}$$

$$K = -x_{m-1} \cdot 2^{m-1} + \sum_{i=0}^{m-2} x_i \cdot 2^i$$

$$1010_{c_2} \rightarrow 2 - 8 = -6$$

$$1001_{c_1}$$

$$1110_{SM}$$

$$-6_{10}$$

$$\boxed{(x_{i-1} - x_i) \cdot 2^i \cdot y}$$

x_i	x_{i-1}	OP
0	0	-
0	1	+4
1	0	-4
1	1	-

$$\text{step } 0: (\cancel{x_{-1}} - x_0) \cdot 2^0 \cdot y +$$

$$\text{step } 1: (\cancel{x_0} - x_1) \cdot 2^1 \cdot y +$$

$$\vdots$$

$$\text{step } i: (\cancel{x_{i-1}} - \cancel{x_i}) \cdot 2^i \cdot y +$$

$$\text{step } m-2: (\cancel{x_{m-3}} - \cancel{x_{m-2}}) \cdot 2^{m-2} \cdot y$$

$$\boxed{P = (-x_{m-1} \cdot 2^{m-1} + \sum_{i=0}^{m-2} x_i \cdot 2^i) \cdot y}$$

$$\begin{aligned} & -x_i \cdot 2^i + x_{i-1} \cdot 2^{i+1} = \\ & = (x_i \cdot 2^i) (-1 + 2) \\ & = x_i \cdot 2^i \end{aligned}$$

Exemplu algoritmul lui Booth:

$$\begin{aligned} X &= -\frac{101}{128} = 1.1100101_{SM} \\ &= 1.00111011 \end{aligned}$$

$$Y = -\frac{15}{32}$$

$$= \underline{1.0111100}_{SM}$$

$$= 1.1000100$$

$$\begin{array}{r} 101 \\ 64 \\ \hline 37 \\ 32 \\ \hline 5 \end{array}$$

28

count	A	Q	M
000	00000000	10011011 0	11000100
	+ 00111100		-M
	—————		
	00111100		
	00011110	01001101 1	
001	00001111	00100110 1	
010	11000100		-M
	11010011		
	11101001	10010011 0	
011 +	00111100		-M
	00100010 1		
	00010010 0	11001001 1	
	+M = 1100 0100		
	-M = 00111100		

100	000001001	011001000 1	
101 +	110001000		-M
	—————		
	11001101		
	11100110	10110010 0	
110	11110011	01011001 0	
111 +	00111100		
	00101111 1	01011000 0	

X _{i+1}	X _i	X _{i-1}	O P
0	0	0	-
0	0	1	+Y
0	(1)	0	+2Y
0	1	1	-2Y
1	0	0	-Y
1	0	1	-Y
1	1	0	-Y
1	1	1	-

$$-Y \cdot 2^i + Y \cdot 2^{i+1} = +Y \cdot 2^i$$

$$+Y \cdot 2^{i+4} - 2Y \cdot 2^i$$

28.02.2023

1 Computer Arithmetic

1.1 Speeding up multiplication Radix-4

x_i	x_{i-1}	OP
0	0	0
0	1	1
1	0	1
1	1	0

$$x = 1000\ 110\ 1:0$$

$$x^* = \bar{1}001\ 0\bar{1}1\bar{1}$$

$$x = 1111\ 1100\ 1:0$$

$$x^* = 0000\ 0\bar{1}00$$

$(x_i)_M$	x_i	x_{i-1}	OP
0	0	0	0
0	0	1	1
0	1	0	1
0	1	1	2
1	0	0	2
1	0	1	1
1	1	0	1
1	1	1	0

$$x = -\frac{101}{128}$$

$$y = -\frac{15}{32}$$

$$M = 1100\ 0100_2 \quad 2M = 11000\ 1000$$

$$-M = 0011\ 1100_2 \quad -2M = 00111\ 1000$$

Count	A	Q	QF-17	M
00	$\begin{array}{r} 0000\ 00000\ 0 \\ + 0001\ 11100 \\ \hline 0001\ 11100 \\ 0000001111 \end{array}$	$\begin{array}{r} 10011\ 011 \\ - \\ 00100\ 110 \end{array}$	$\begin{array}{r} -1 \\ 0 \\ -1 \end{array}$	11000100
01	$\begin{array}{r} 0000111100 \\ + 001001011 \\ \hline 000010010 \end{array}$	$\begin{array}{r} 11001001 \\ - \\ 10110010 \end{array}$	$\begin{array}{r} 1 \\ -1 \\ 2M \end{array}$	
10	$\begin{array}{r} 11000\ 1000 \\ + 110011010 \\ \hline 111100110 \end{array}$	$\begin{array}{r} 10110010 \\ - \\ 10110010 \end{array}$	$\begin{array}{r} 0 \\ 0 \end{array}$	
11	$\begin{array}{r} 001111000 \\ + 001011110 \\ \hline 000100010 \end{array}$	$\begin{array}{r} 101100 \\ - \\ 101100 \end{array}$	$\begin{array}{r} 0 \\ 0 \end{array}$	

No transformation to final

$$x = 1001\ 10111\ 0$$

$$1001\ 10111\ 0$$

$$x^* = \bar{1}010\ \bar{1}101 \quad \text{Booth codificat}$$

$$x^{**} = \bar{1}\ 2\ \bar{1}\ \bar{1}$$

Section 2

1.2. Redundant digit set representation

Radix 2: $\{0, 1\} \rightarrow \{\bar{1}, 0, 1\}$

Radix 4: $\{0, 1, 2, 3\} \rightarrow \{\bar{2}, \bar{1}, 0, 1, 2\}$

$$10\bar{1}1 = 1 \cdot 2^0 + 0 \cdot 2^1 + 1 \cdot 2^2 + 1 \cdot 2^3 = 7$$

0 { 0

1 { 1
0

$$\begin{array}{r} 1001 \\ 0010 \\ \hline 0111 \end{array}$$

\bar{1} { 0

$^{+3} 2^{-1} 0$

$$102\bar{2}\bar{1} = 1 \cdot 4^4 + 2 \cdot 4^2 - 2 \cdot 4^1 - 1 \cdot 4^0 \quad - \text{Radix 4}$$
$$= 249$$

0 → 00

$$\begin{array}{r} 01 00 10 00 00 - \\ 00 00 00 10 01 \\ \hline 01 00 01 01 11 \end{array}$$

$$\begin{array}{r} 23 + \\ 256 \\ \hline 279 \end{array}$$

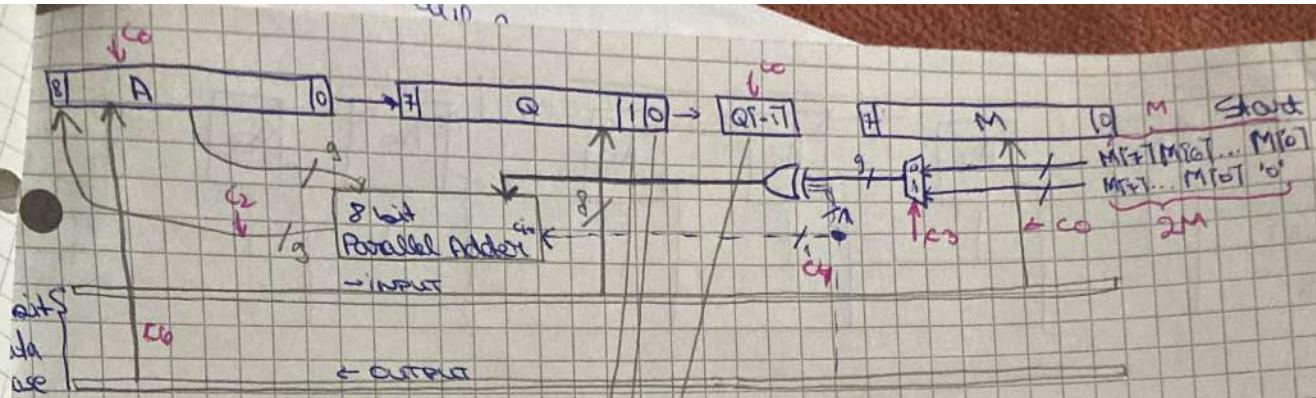
1 → 01
00

2 → 10
00

$\bar{1}$ → 00

$\bar{2}$ → 00

1.3. Hardware implementation of radix 4 - Booth

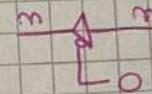


\leftarrow COUNT

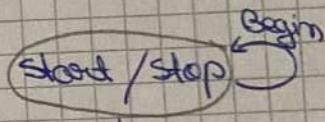
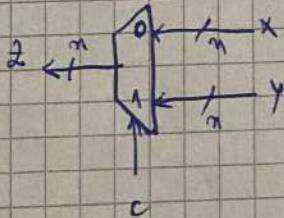
COUNT₃

Begin
Clock
Reset
End

Control
Unit

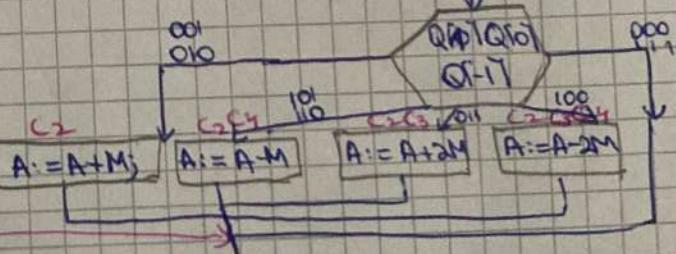


000	-0
0·01·2	1
010·1	
011·2	
100·1	
101·2	
110	
111	0



$A := 0, Q[7:-1] := 0,$
 $COUNT := 0, M := INBUS$

$Q[7:-1] := INBUS$



CONTINUARE APRODU :

Count 3

$A[6:-1], Q[7:-1] := A[6:-1], Q[7:-1]$
 $A[7] := A[8], A[8] := A[8]$
 $COUNT := COUNT + 1;$

I 6

M

0

GND

M[0]

↓

0

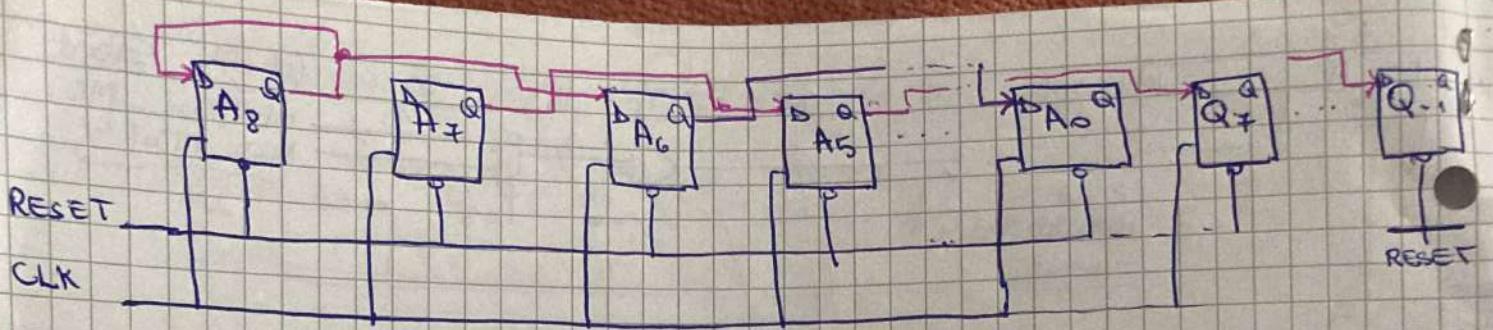
OUTBUS = A[8:-1]

C5

OUTBUS = A[0].Q[7:-1] → Start / Stop

C6

End



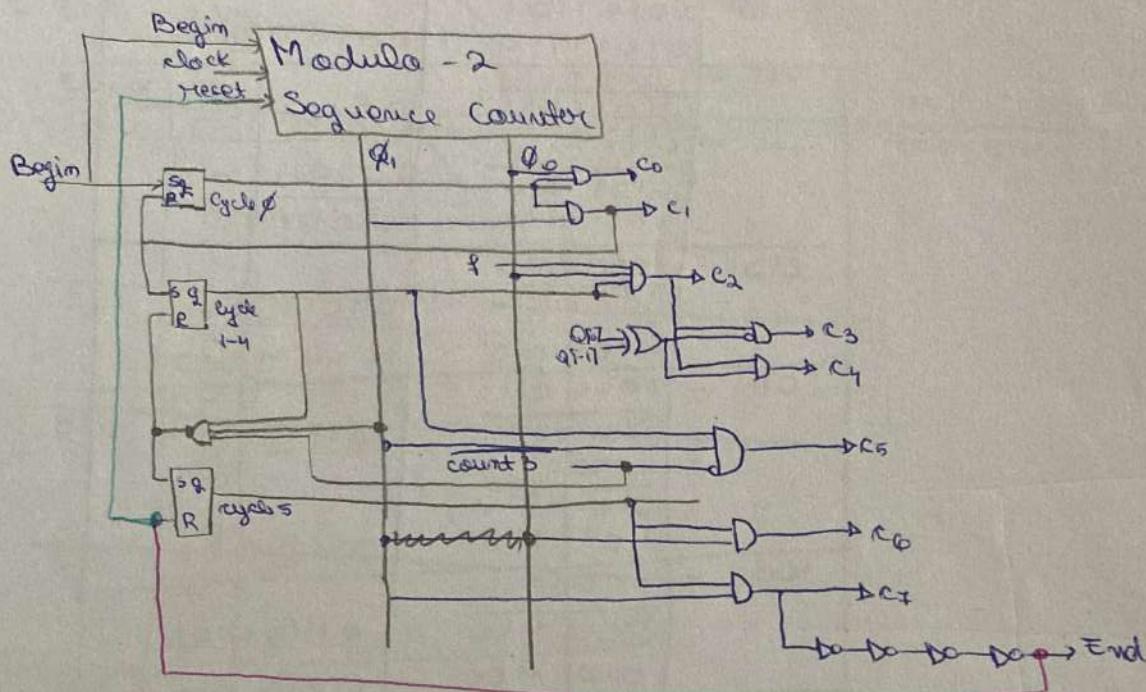
Cours 3

4.03.2023

Episodeul 2

Sequence Counter

CU synthesis



§.3 Division

$$\begin{array}{r}
 \overbrace{\quad\quad\quad}^{\text{Dividend}} \\
 5771 \\
 \hline
 \overbrace{\quad\quad\quad}^{\text{Divisor}} \quad | \quad \overbrace{135}^{\text{Quotient}} \\
 540 \\
 \hline
 = 371 \\
 \overbrace{\quad\quad\quad}^{\text{Reminder}} \quad | \quad \overbrace{270}^{\text{Quotient}} \\
 101
 \end{array}$$

Binary

$$\begin{array}{r}
 1010\ 0010 \\
 1001 \\
 \hline
 0001011 \\
 -1001 \\
 \hline
 00100
 \end{array}$$

$$\begin{array}{r} 38 \\ \times 128 \\ \hline 1656 \end{array}$$

$$\begin{array}{r} 166 \\ \times 9 \\ \hline 1494 \end{array}$$

size of the quotient = size of the divisor = m bits

$$\Rightarrow \text{Dividend size} = 24 - 1 \text{ bits}$$

1.3.1.

$$\leftarrow 2x_i - q_i M$$

Restoring division

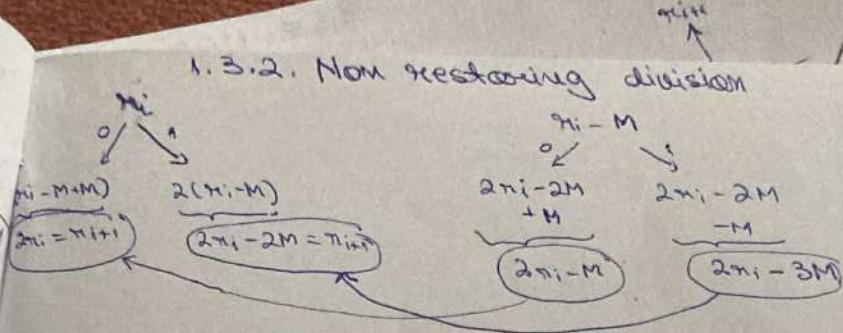
$$1011010001011 = 577_{10}$$

$$135_{10} = 10000111$$

COUNT	A	Q	M
000	00100101 - 10000111 ----- 01100110	00010110 00010110	10000111
-	10000111	00101110	
+ 100	10000111 - 01010011 ----- 01011010	00101110	
-	10000111 01011010 ----- 10110100	01011000	
+ 010	10000111 00101110 01011010	01011000	
-	10000111 01011010 ----- 10110101	10110010	
+ 011	10000111 01011010 10110101	10110010	
-	10000111 01011010 ----- 10110101	01100100	
+ 100	10000111 00101110 01011010	01100100	
-	10000111 01011010 ----- 10110101	11001010	
+ 101	10000111 01010101 10000111	11001010	
-	10000111 01010101 ----- 01011100	10101010	
+ 110	10000111 00110010 01100101	10101010	
-	10000111 00110010 ----- 01100101	10101010	
+ 111	10000111 01011110 10000111	10101010	
-	10000111 01011110 ----- 01100101	42	
			101_{10}

Example: $n=8 \Rightarrow 13$ counts. op.

1.3.2. Non restoring division



$$\begin{array}{r} 64 \\ + 37 \\ \hline 101 \end{array}$$

1.3.3 Division

COUNT	S	A	Q	M
000	0	0010 1101	0001 0110	1000 0111
	-	1000 0111		
	①	1010 0110		
	1			
001	+			
	①			
	1			
010	+0			
	①			
	0			
011	-			
	①			
	1			
100	+			
	①			
	0			
101	-			
	①			
	1			
110	+			
	①			
	0			
111	-			
	①			
	+			

SRT Division

14.03.2023

Ch 4

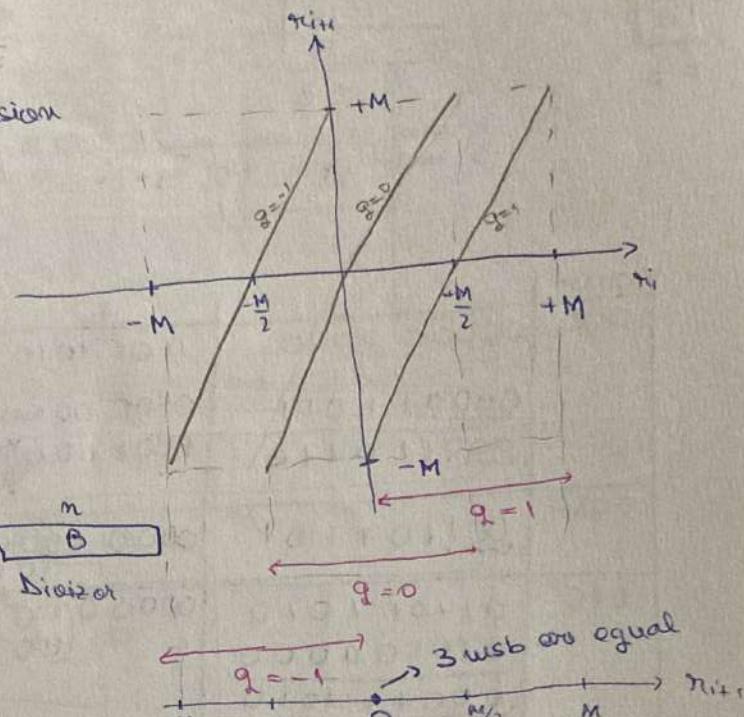
$$x_{i+1} \leftarrow 2x_i - q_i M$$

division

$$|x_i| < M$$

$$|x_{i+1}| < M$$

$$q_i \in \{-1, 0, 1\}$$



1. Shifting over zeros Shift P.A.B \leftarrow $\frac{m}{2}$ positions (k leading 0s)

2. For $i=0$ to $m-1$

a) if 3 msb $\overset{\text{of } P}{\sim}$ are equal
 $q_i = 0$

Shift P.A one position left

b) if 3 msb are not equal and P positive
 $q_i = 1$

Shift P.A one position left

$P \leftarrow P - B$

c) if 3 msb are not equal and P negative
 $q_i = -1$

Shift P.A. one position left

$P \leftarrow P + B$

3. Correction

if $P < 0$

$P \leftarrow P + B$

$Q \leftarrow Q - 1$

4. Shift P right k positions

$$\begin{array}{r}
 213 \\
 20 \\
 \hline
 = 13 \\
 10 \\
 \hline
 = 3
 \end{array}
 \quad \left| \begin{array}{c} 5 \\ 42 \end{array} \right.$$

$$\begin{array}{r}
 218 \\
 20 \\
 \hline
 = 18 \\
 15 \\
 \hline
 = 3
 \end{array}
 \quad \left| \begin{array}{c} 5 \\ 43 \end{array} \right. - x$$

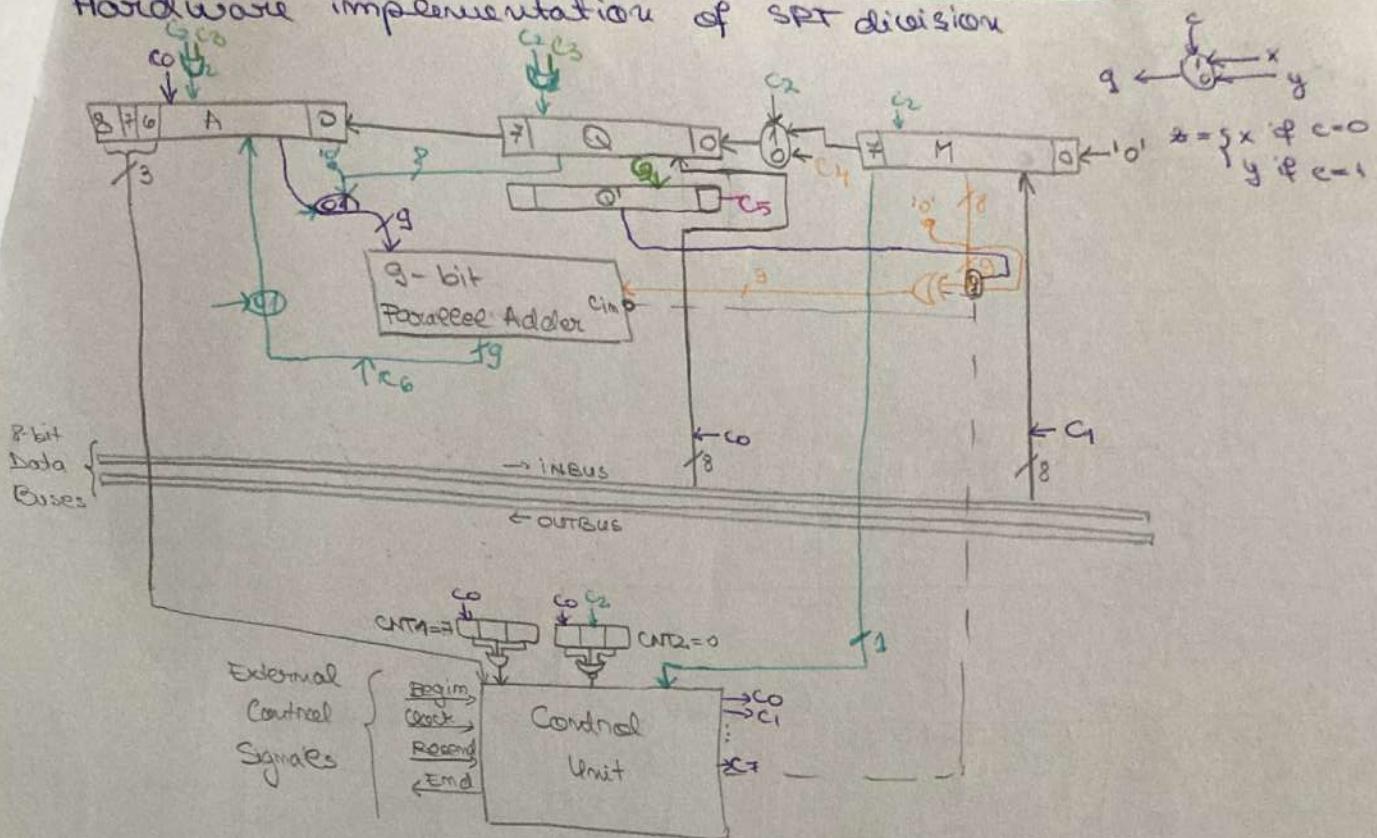
$$\begin{aligned}
 218_{10} &= 1101 \ 1010 = A \\
 45_{10} &= 0010 \ 1011 \\
 5_{10} &= 0000 \ 0101 = B \\
 3_{10} &= 0000 \ 0011
 \end{aligned}$$

3 msb of P are equal $\Rightarrow q_i = 0$

3 msb of P NOT equal, $P > 0 \Rightarrow q_i = 1$, Shift si $P = P - B$

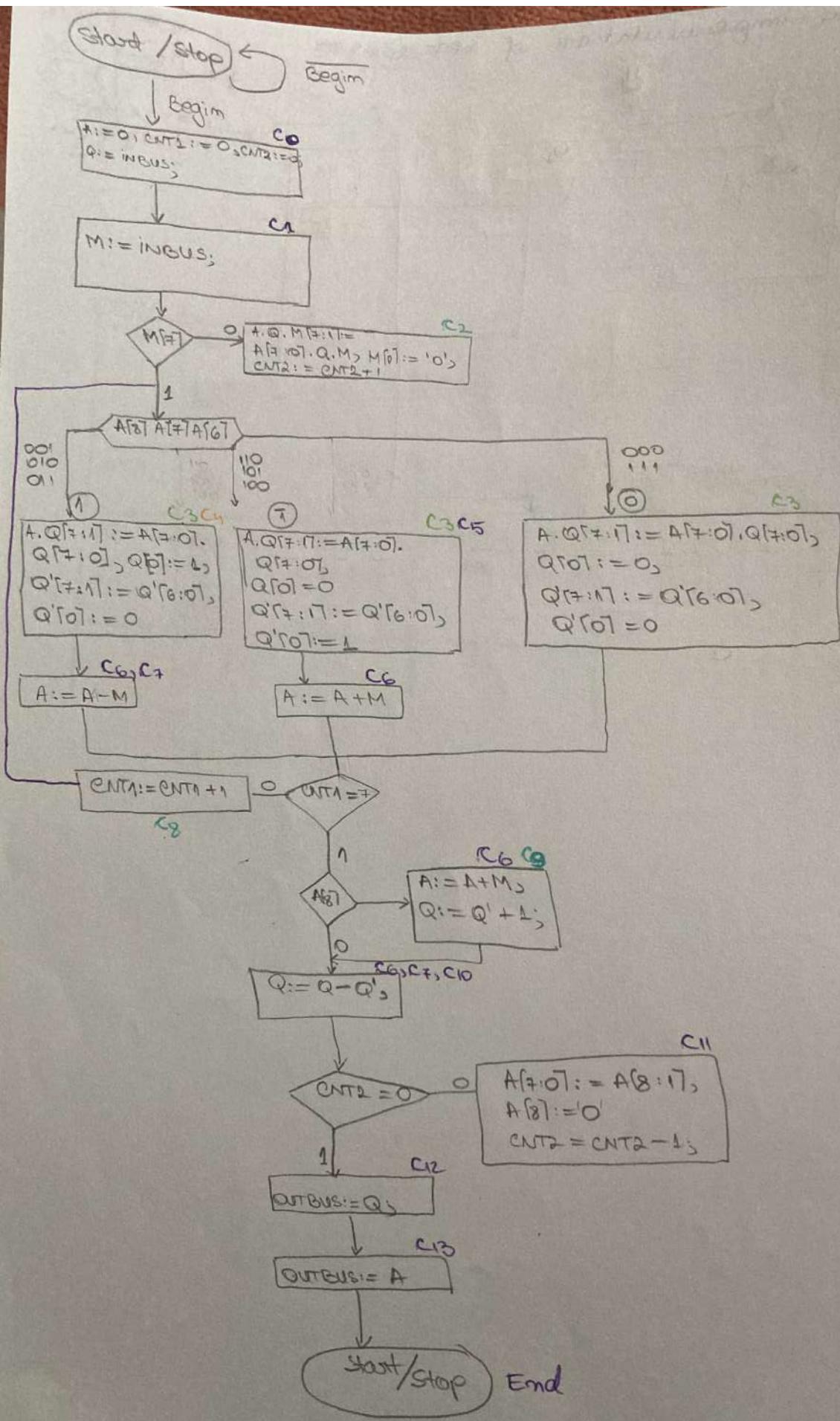
COUNT	P	A	B	
000	00000 0000	1101 1010	0000 0101 1010 0000	$K = 5$ $\ll 5$ $\ll 1$
001	00110 1101	0000 0000		$\ll 1$
010	01101 1010 - 01010 0000	0000 0000 0000 0001		$0 \rightarrow 0 \ A$ $0 \rightarrow 0 \ A'$ $1 \rightarrow 1 \ A$ $1 \rightarrow 1 \ A'$
	00011 1010			T → 0 A 1 A'

Hardware implementation of SRT division



$$x = \begin{cases} x & \text{if } c=0 \\ y & \text{if } c=1 \end{cases}$$

$$q \leftarrow \begin{cases} 1 & \text{if } x \\ 0 & \text{if } y \end{cases}$$





HOMANIT

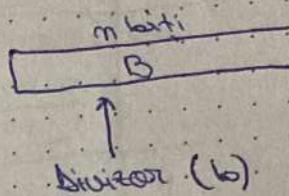
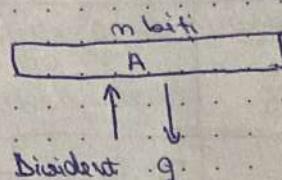
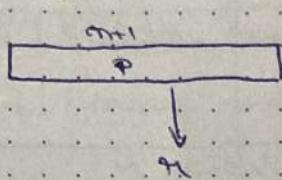
21.03.2023

Cours 5

1.4.5 Speeding up SRT Division

SRT Radix-4

$$q_i \in \{2, 1, 0, -1, -2\}$$



$$|x_i| < b \Rightarrow |x_{i+1}| < b$$

$$x_{i+1} \leftarrow 4x_i - q_i b$$

$$x_i = b$$

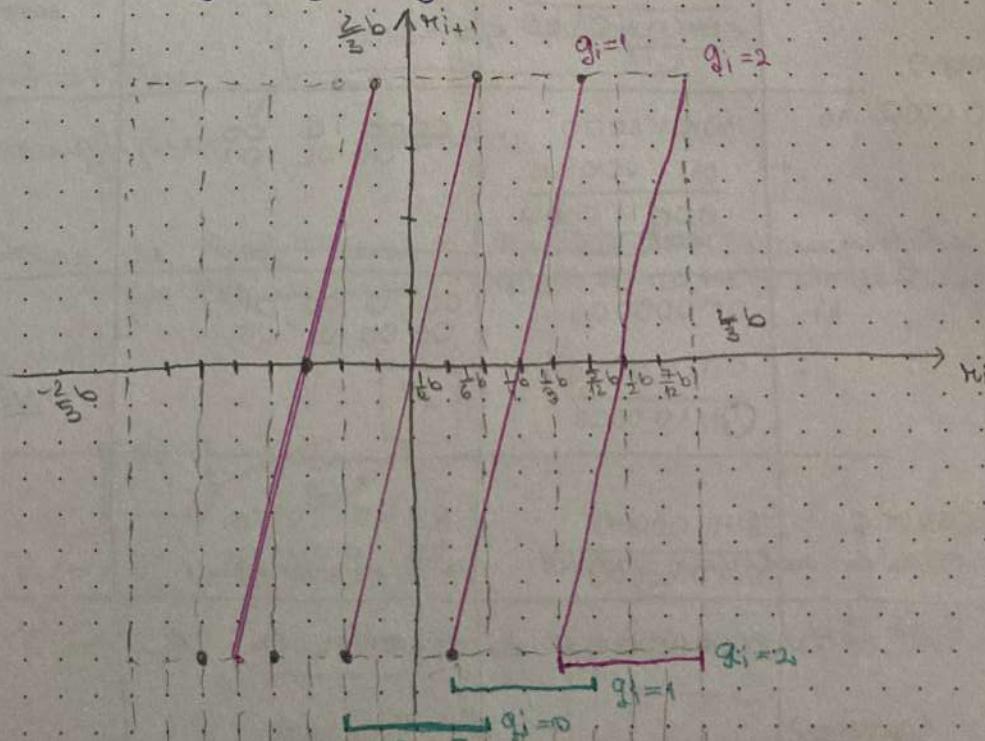
$$x_{i+1} = 4b - q_i b$$

Dacă am pure $|x_i| < \frac{2}{3} b$

$$x_i = \frac{2}{3} b$$

$$x_{i+1} = \frac{8}{3} b - q_i b \quad q_i = 2$$

$$\Rightarrow x_{i+1} = \frac{8}{3} b - \frac{6}{3} b = \frac{2}{3} b$$



$$\begin{array}{r}
 211 \\
 18 \\
 \hline
 = 31 \\
 30 \\
 \hline
 = 1
 \end{array}$$

Shifter: deelbaar

COUNT	P	A	B
00	00000 0000 0 0001 1001 0011 01, 001	1101 0011 0110 0000 1000 00 00 00	00000 0110 1100 0000 b=12
01	1101 00110 1100 00000 0001 00110	00 00 00 10 00 00	
1 → 01 00	-		
2 → 10 00	10	0100 11 000 0110 00000 1110 11 000	00 00 10 01 00 00 00
7 → 00 01	11	1011 00000 0110 00000 0001 00000	00 10 01 00 00 00 00 01
I → 00 10	+ CONVERSIE SHIFTERE	0000 0 0001 Reminder	00 1000 011 Quotient

$$\begin{array}{r}
 202 \\
 14 \\
 \hline
 = 62 \\
 56 \\
 \hline
 = 6
 \end{array}$$

$$B = 0111 00000$$

$$QB = 1100 0000$$

CORECTIE +

SHFT 13

CONVERSIE

Reminder

Quotient



Cap. 2.

Performante sistemelor de calcul

→ 3 tipuri de sisteme de calcul

1. Desktop computers 40 years of progress

2. Servers and datacenters

- scientific

- many users at the same time

3. Embedded computers

Performanta

- Time

- Space (Memory, circuits)

- Power consumption

- Dependability

Runtine

Execution Time

Throughput

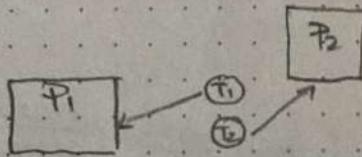
} Desktop computers

- Servers

sisteme de Real Time - Embedded computers

- performance : worst case : deadline

Exemplu:



- ① $P_1 \rightarrow P_1$, better runtime \Rightarrow better runtine \Rightarrow better throughput
- ② $P_1 \rightarrow P_1, P_2 \Rightarrow$ better throughput \Rightarrow better runtine

$$\text{Performance}_x = \frac{1}{\text{Execution time}_x}$$

Example:

Computer A runs a program in 15 seconds

Computer B runs the same program in 10 seconds

$$\frac{\text{Perf. B}}{\text{Perf. A}} = \frac{\text{Exec. time A}}{\text{Exec. time B}} = \frac{15}{10} = 1,5$$

Benchmark units



Cours 6

2.2 Performance metrics

$$\text{Performance}_x = \frac{1}{\text{Execution time}_x}$$

- Real programmes (run)
- Not
 - toy programmes
 - synthetic
 - kernels
- transitivity $A > B \wedge B > C$
 $\Rightarrow A > C$

Dec	Binary	
10^3	2^{10}	1K
10^6	2^{20}	1 Mi
10^9	2^{30}	1 Gi
10^{12}	2^{40}	1 Ti
10^{15}	2^{50}	1 Pi peta
10^{18}	2^{60}	1 Ei exa
10^{21}	2^{70}	1 Zi zetta
10^{24}	2^{80}	1 Yi yotta

Benchmark suites

SPEC : Standard Performance Evaluation Corporations
 1989 , 92 , 95 , 2000 , 2006 , 2014

- CPU SPEC
- SPEC graphics
- SPEC Java

2.2.1 Reporting performance results

SPEC Ratio

$$\frac{\text{SPEC Ratio A}}{\text{SPEC Ratio B}} = \frac{\frac{\text{Execution time reference}}{\text{Execution time A}}}{\frac{\text{Execution time reference}}{\text{Execution time B}}} = \frac{\text{Execution time B}}{\text{Execution time A}} \rightarrow \frac{\text{Performance A}}{\text{Performance B}}$$

$$= \frac{\text{Execution time B}}{\text{Execution time A}} \rightarrow \frac{\text{Performance A}}{\text{Performance B}}$$

$$\text{Geometric mean} = \sqrt[n]{\prod_{i=1}^n \text{Sample}_i}$$

SPEC ratio \downarrow
program

1. The geometric means of ratios = the ratio of geometric means
2. The ratio of geometric means

$$\frac{\text{Geometric Mean A}}{\text{Geometric Mean B}} = \frac{\sqrt[n]{\prod_{i=1}^n \text{Spec Ratio } A_i}}{\sqrt[n]{\prod_{i=1}^n \text{SPEC ratio } B_i}} = \boxed{n=29}$$

$$= \sqrt[n]{\prod_{i=1}^n \frac{\text{SPEC Ratio } A_i}{\text{SPEC Ratio } B_i}} = \sqrt[n]{\prod_{i=1}^n \frac{\frac{\text{Execution time ref.}}{\text{Execution time } A_i}}{\frac{\text{Execution time ref.}}{\text{Execution time } B_i}}}$$

$$= \sqrt[n]{\prod_{i=1}^n \frac{\text{Execution time } B_i}{\text{Execution time } A_i}} = \sqrt[n]{\prod_{i=1}^n \frac{\text{Performance } A_i}{\text{Performance } B_i}}$$

$$\text{stddev} = \sqrt{\sum_{i=1}^n (\text{Sample}_i - \text{Mean})^2}$$

\uparrow
SPEC ratio $_i$

Geometric standard deviation

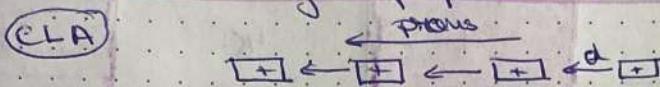
$$\text{Geometric mean} = \exp \left(\frac{1}{n} \cdot \sum_{i=1}^n \ln(\text{Sample}_i) \right)$$

$$\text{gstdev} = \exp \left(\sqrt{\frac{1}{n} \cdot \sum_{i=1}^n [\ln(\text{Sample}_i) - \ln(\text{Geom. mean})]^2} \right)$$

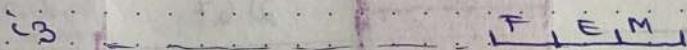
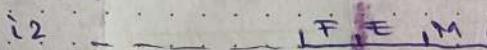
2.3. Quantitative principles of computer design

SLD - System level design

- L. Take advantage of parallelism (when possible)

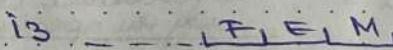
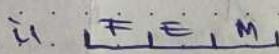


Instruction-level parallelism



19

→ time



→ time

2. Principle of locality

Obs: Programms tend to reuse data and instructions they have used recently.

Rule of

A programme spends 90% of the time on 10% of the code

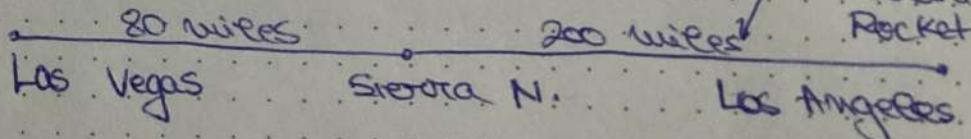
- Mostly applied on instructions

- Also works on data

- 2 types of locality \leftrightarrow temporal locality
 \rightarrow spatial locality

- ### 3. Focus on the common case

Arndtke's law



HOMANIT GMBH & CO. KG
Bahnhofstraße 30
37412 Herzberg am Harz

VERTRIEB · SALES DEPARTMENT · SERVICE COMMERCIAL
Tel. +49 5521 84-0 · Fax +49 5521 84-269
info@homanit.de · www.homanit.de

SPRZEDAZY
Tel. +48 94 3100571
info@homanit.pl · www.homanit.pl

THAMON

Vehicle for the second part

Vehicle for the 2nd part	Hours for the 2nd part	Speed up in the desert	Hours for the entire trip	Speedup for the entire trip
Feet	50	1	40	1
Bike	20	2.5	70	1.75
Hyundai	4	12.5	24	2.9
Ferrari	1.67	30	21.67	3.23
Rocket car	1	50	21	3.3

Speedup for the entire system =

Performance of the task using the speedup where possible
Performance of the task without the speedup

$$= \frac{\text{Execution time of task without speedup}}{\text{Execution time of task using speedup where possible}} =$$

$$= \frac{\text{Execution time of task without speedup}}{\text{Exec. time of task without speedup} \left[(1 - \text{Fraction enhanced}) + \frac{\text{Fraction enh.}}{\text{Speedup}} \right]}$$

$$= \frac{1}{(1 - \text{Fraction enhanced}) + \frac{\text{Fraction enhanced}}{\text{Speedup}}} \rightarrow \frac{200}{280}$$

Speedup in the desert



Limit of system enhancement

$$\text{line speedup} \rightarrow \text{Speedup for the entire system} = \frac{1}{1 - \text{Fraction enhanced}}$$

2.4 Computer performance equation

$$\text{Clock cycle time} = \frac{1}{\text{Clock rate}}$$

$$\text{Clock cycles for a program} = \underbrace{\text{Instruction count (IC)}}_{\substack{\uparrow \\ \text{Integer}}} \times \text{Clock cycles per Instruction (CPI)}$$

$$\text{CPU time} = \text{IC} \times \text{CPI} \cdot \text{Clock cycle time}$$

$$\begin{aligned} \text{CPU time} &= \frac{\cancel{\text{Instructions}}}{\text{Programme}} \cdot \frac{\cancel{\text{Clock cycles}}}{\text{Instructions}} \cdot \frac{\cancel{\text{Seconds}}}{\text{Clock cycles}} = \\ &= \frac{\text{seconds}}{\text{programme}} \end{aligned}$$

Improvement = reduce IC

RISC processors

CISC

4.04.2023

Cours 7.

= continue =

Example 1

- Frequency of FP operations = 25 %
- $CPI_{FP} = 4$ clock cycles
- CPI of other instructions = 1,33 clock cycles
- Frequency of FPSQRT = 2 %
- $CPI_{FPSQRT} = 20$ clock cycles

Design optimization alternative

- a) decrease CPI_{FPSQRT} to 2 clock cycles
- b) decrease CPI of all FP to 2.5 clock cycles

$$CPU \text{ time reference} = IC_n \cdot CPI_n \cdot \text{clock cycle time}_n$$

$$CPI_n = 0,25 \cdot 4 + 0,75 \cdot 1,33 = 2 \text{ c.c}$$

Alternative A

$$CPU \text{ time A} = IC_n \cdot CPI_A \cdot \text{clock cycle time}_A$$

$$CPI_A = 2 - 18 \cdot 0,02 = 1,64 \text{ c.c}$$

Amdahl's Law

Alternative B

$$CPU \text{ time B} = IC_n \cdot CPI_B \cdot \text{clock cycle time}_B$$

$$CPI_B = 0,25 \cdot 2,5 + 0,75 \cdot 1,33 = 1,625 \text{ c.c} \quad \checkmark$$

Wessel 8

Exemplul 2

CPU A RISC

20 % are branch instruction
 \Rightarrow 2 % compare instruction

CPI branch = 2 cc

CPI all other = 1 cc

$a \rightarrow r_1$
 $b \rightarrow r_2$
 CMP r_1, r_2
 BNE address

if $a = b$ then ~~end~~
 jump address r_1
 else \dots

$\dots r_1$ ~~address~~ $\dots r_2$

B address \dots

address $\dots r_2$

address 2 \dots

CPU B CISC

BNE $r_1, r_2, \text{address}$

Clock cycle time is extended with 25 %

$$\text{CPU time A} = IC_A \cdot CPI_A \cdot \text{clock cycle time A} = IC_A \cdot 1,2 \cdot \text{clock cycle time A}$$

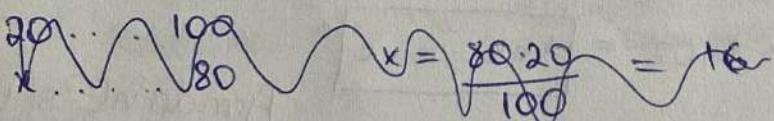
de degradația clock-ului \Rightarrow waste time

$$CPI_A = 0,2 \cdot 2 + 0,8 \cdot 1 = 1,2 \text{ cc}$$

$$CPI_B = IC_B \cdot CPI_B \cdot \text{clock cycle time B}$$

$$IC_B = (100 - 20) \% IC_A$$

$$\text{Clock cycle time B} = 1,25 \cdot \text{clock cycle time A}$$



$$20 \dots 80 \quad x = \frac{80 - 100}{100} = \frac{1}{4} \quad x = 25 \%$$

$$x \dots 100$$

$$CPI_B = 0,25 \cdot 2 + 0,75 \cdot 1 \\ = 1,25 \text{ cc}$$

$$\text{CPU time B} = IC_A \cdot 0,8 + 1,25 \cdot 1,25 \cdot \text{clock cycle time A}$$

$$= IC_A \cdot 1,25 \cdot \text{clock cycle time A}$$

Concluzie: Sometimes less is more ☺!

Example 3

Load-store machine

$LDR \quad r_3, [r_4]$
 $ADD \quad r_2, r_2, r_3$
 $STR \quad r_2, [r_5]$

Register memory

$ADD \quad r_2, r_3, [r_4]$

instruction

Instruction	CPI	Frequency
ALU	1	43 %
Load	2	21 %
Store	2	12 %
Branch	2	24 %

25% of ALU instructions have an operand taken from memory and placed in a register

New machine

- a new type of ALU instructions
- ALU CPI = 2 CC
- Branch CPI = 3 CC

Relevance:

$$\begin{aligned} \text{CPU time old} &= IC_{old} \cdot CPI_{old} \cdot \text{Clock cycle time old} \\ &= IC_{old} \cdot 1,57 \cdot \text{Clock cycle time old} \quad \checkmark \end{aligned}$$

$$IC_{old} = 1 \cdot 0,43 + 2 \cdot 0,57 = 1,57$$

$$\text{CPU time new} =$$

$$IC_{new} = IC_{old} (1 - 0,25 \cdot 0,43)$$

$$\text{Clock cycle time old} = \text{Clock cycle time new}$$

$$\begin{aligned} &\cancel{(0,43 \cdot 0,75)} \cdot 1 + \cancel{(0,43 \cdot 0,25)} \cdot 2 + \cancel{(0,21 - 0,43 \cdot 0,25)} \cdot \\ &\cancel{+ 0,12 \cdot 2 + 0,24 \cdot 3} \end{aligned}$$

$$CPI_{new} =$$

$$CPI_{new} = \frac{(0,43 \cdot 0,75) \cdot 1 + (0,43 \cdot 0,25) \cdot 2 + (0,21 - 0,43 \cdot 0,25) \cdot 2 + 0,12 \cdot 2 + 0,24 \cdot 3}{1 - 0,43 \cdot 0,25}$$

$$\begin{aligned} \text{CPU time new} &= IC_{old} \cdot (1 - 0,25 \cdot 0,43) \cdot \frac{\cancel{(1 - 0,43 \cdot 0,25)}}{\cancel{(1 - 0,43 \cdot 0,25)}} \cdot \text{Clock cycle time old} \\ &= IC_{old} \cdot 1,908 \cdot \text{Clock cycle time old} \quad \times \end{aligned}$$

2.5. Other metrics

MIPS = Million Instructions Per Second

$$\text{MIPS} = \frac{IC}{\text{CPU time} \cdot 10^6} = \frac{f_c}{f_c \cdot CPI \cdot \text{clock cycle time} \cdot 10^6} = \\ = \frac{\text{Clock rate}}{CPI \cdot 10^6}$$

Problems with MIPS

1. Cannot compare 2 components with distinct instruction sets
2. MIPS varies for the same computer for different programs
3. Sometimes MIPS is misleading.

Exemplu 4 The same L/S machine from Exemplu 3

Instructions	Frequency	CPI
ALU	43 %	1
Load	21 %	2
Store	12 %	2
Branch	24 %	2

Clock cycle time = 20 ns

ALU ops reduced with 50% by a compiler optimization

$$\text{Frequency} = \frac{1}{20 \text{ ns}} = \frac{1}{20 \cdot 10^{-9} \text{ s}} = 50 \text{ MHz}$$

$$\text{CPU time unoptimized} = IC_u \times CPI_u \times 20 \text{ ns} = IC_u \times 1,57 \times 20 \text{ ns} \\ = IC_u \times 31,4 \text{ ms}$$

$$\text{MIPS}_u = \frac{50 \cdot 10^6 \text{ Hz}}{1,57 \cdot 10^6} = \frac{50}{1,57} = \boxed{31,85}$$

18.04.2023

Werkel 8

Example

Instruction type	Frequency	CPI
ALU	43%	1
LOAD	21%	2
STORE	12%	2
BRANCH	24%	2

$$\text{Clock cycle time} = 20 \text{ ns}$$

$$\text{Clock frequency} = \frac{1}{20 \text{ ns}} = 50 \text{ Hz}$$

$$\text{CPU time original} = IC \cdot CPI_{old} \cdot 20 \text{ ns} = IC_0 \cdot 1,57 \cdot 20 = 3,14 IC_0$$

$$CPI_{old} = 0,43 \cdot 1 + 0,57 \cdot 2 = 1,57$$

$$\text{MIPS} = \frac{50 \cdot 10^6}{1,57 \cdot 10^6} = 31,85$$

$$IC_{new} = IC_0 \left(1 - \frac{0,43}{2}\right) \Rightarrow CPI_{new} = \frac{\frac{0,43}{2} \cdot 1 + 0,57 \cdot 2}{\left(1 - \frac{0,43}{2}\right)} = 1,73$$

$$\text{CPU new} = IC_0 \left(1 - \frac{0,44}{2}\right) \frac{\frac{0,43}{2} \cdot 1 + 0,57 \cdot 2}{\left(1 - \frac{0,43}{2}\right)} \cdot 20 \text{ ns} = IC_0 \cdot 2,71 \text{ ms}$$

(Mai bunt)

$$\text{MIPS new} = \frac{50 \cdot 10^6}{1,73 \cdot 10^6} = 28,9$$

$$\text{Relative MIPS} = \frac{\text{Execution time reference machine}}{\text{Execution time x}} \cdot \text{MIPS ref. machine}$$

$$\text{MFLOPS} = \frac{IC_{fp}}{\text{CPU time} \cdot 10^6}$$

Lineworm Loops

Floop	Norm Floop
+/-/·	1
Sqrt	2
Exp Sine	8

Curs 9

3. Memory hierarchy

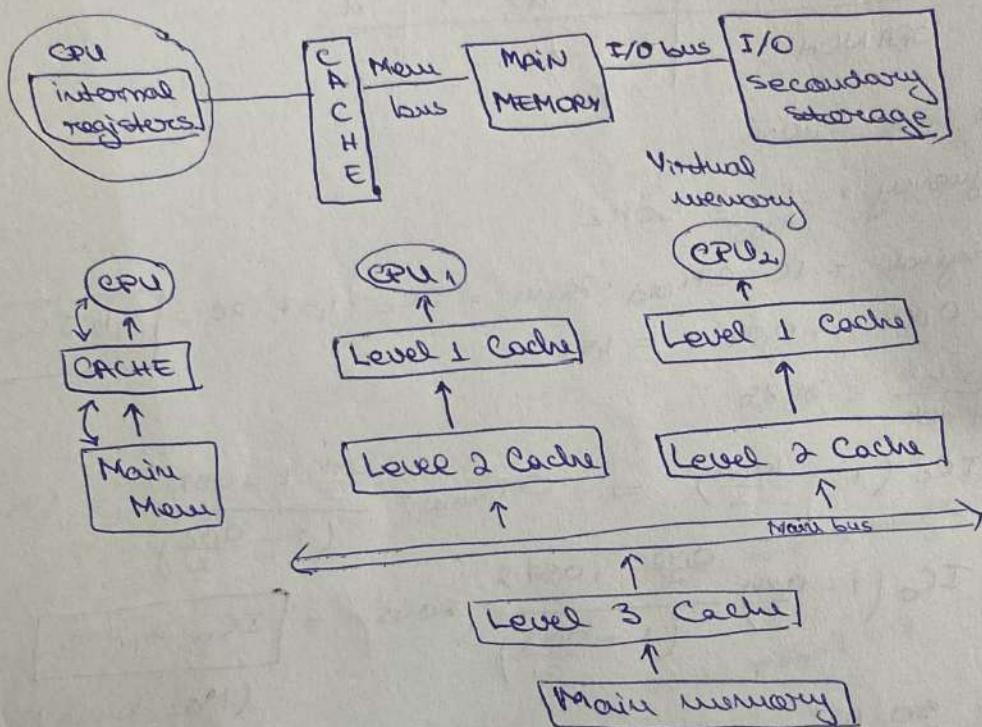
Capitolul III : Hierarhia de memorie

3.1. Introduction

3.2. Cache mapping

Locality Principle

Make the common case faster



Taxonomy & basic notions

Addressable unit → generally $1B = 1B$ Byte

1 word = 2^n bytes

1 block = 2^p words

1 cache hit

1 cache miss → Miss Rate $\approx 10\%$.

Cache coherency

Cache - unified

- split → Data
Instruction

3.2.1. Direct mapping

Ans.

2^m blocks in M_2 (Main Memory)

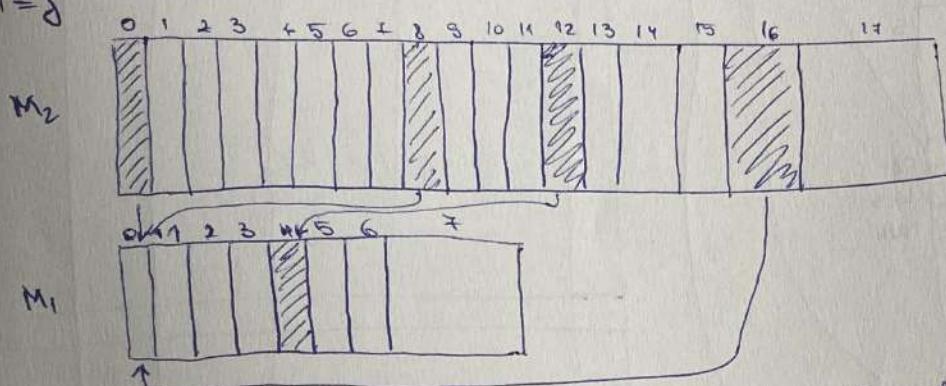
2^n blocks in M_1 (Cache)

$m < n$

j indexes the 'blocks' in the main memory

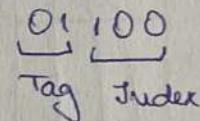
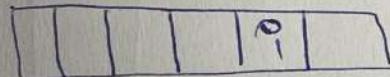
i indexes the blocks in the cache M_1

$i = j$



$$12 \bmod 8 = 4$$

Tag



Address word:

Tag	Index	Word offset	Byte offset
$\langle m-n \rangle$	$\langle n \rangle$	$\langle p \rangle$	$\langle \pi \rangle$

Example 1:

Memory principle - 2²⁴ bits

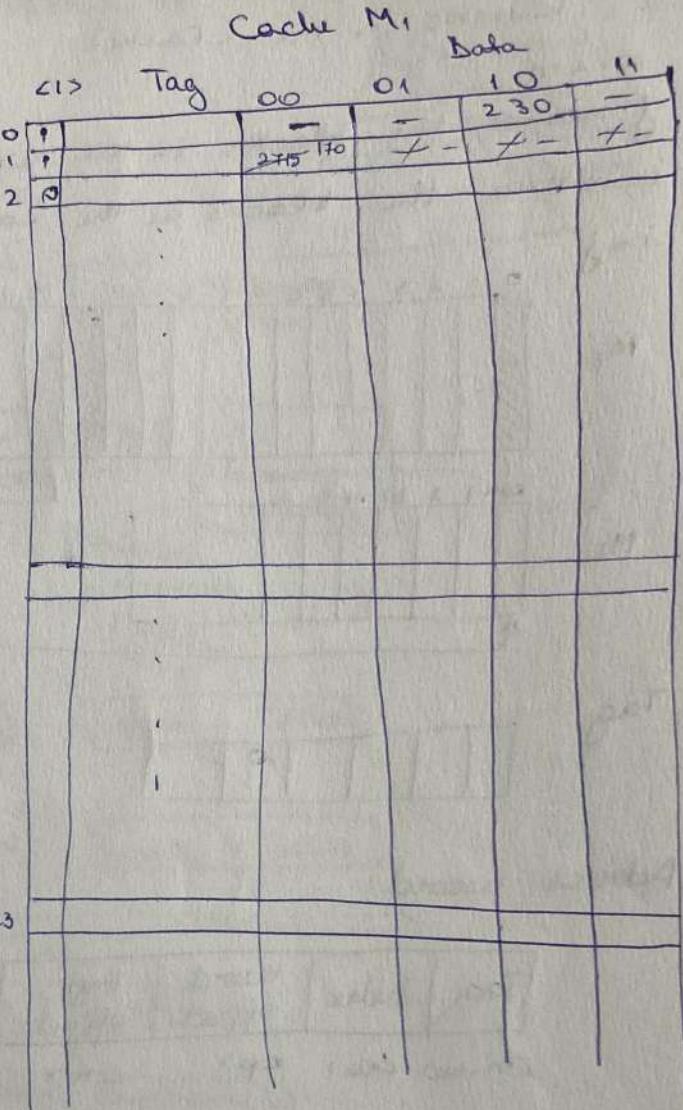
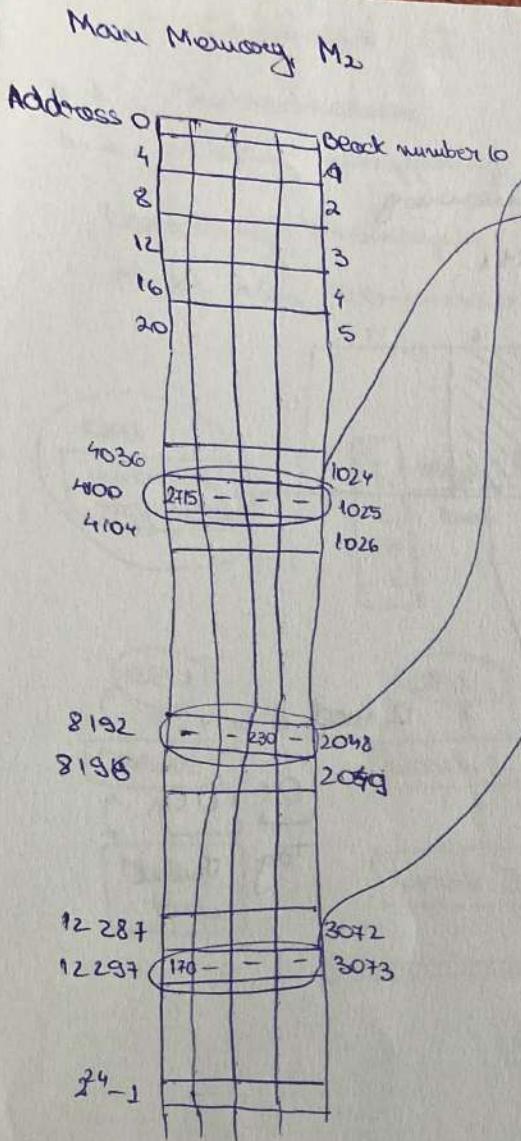
1 byte = 1 byte

1 block = 4 words

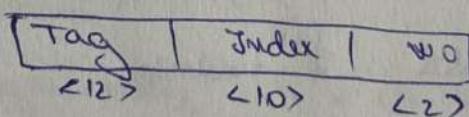
Cache data dimension = 1 KB

Program:

Address	Code
8194	230
4100	2715
12292	170



Address word



	Tag	Index	WD
9196	000000 000010	0000000000	10
4100	000000000001	00 0000000001	00
12292	00 0000000011	00000000001	00

Lecture 9

3. Memory hierarchy

3.1. Introduction

3.2. Mapping

3.2.1. Direct mapping

Example : Intrinsicity Fast MATH

- MIPS architecture
- 32 bit words
- split (16 kB)
- 1 block = 16 words
- byte addresses

Address word

$$1 \text{ word} = 32 \text{ bits} = 4 \text{ bytes} = 2^2 B$$

$$1 \text{ block} = 16 \text{ words} = 2^4 \text{ words}$$

$$\text{Size } (M_2) = 2^4 \cdot 2^2 B = 2^6 B$$

$$= 2^{32} B = 2^2 \cdot 2^{30} B$$

$$= 4 GB$$

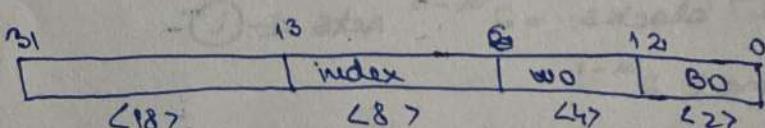
$$= \frac{2^{32} B}{2^6 B} = 2^{26} B \text{ blocks}$$

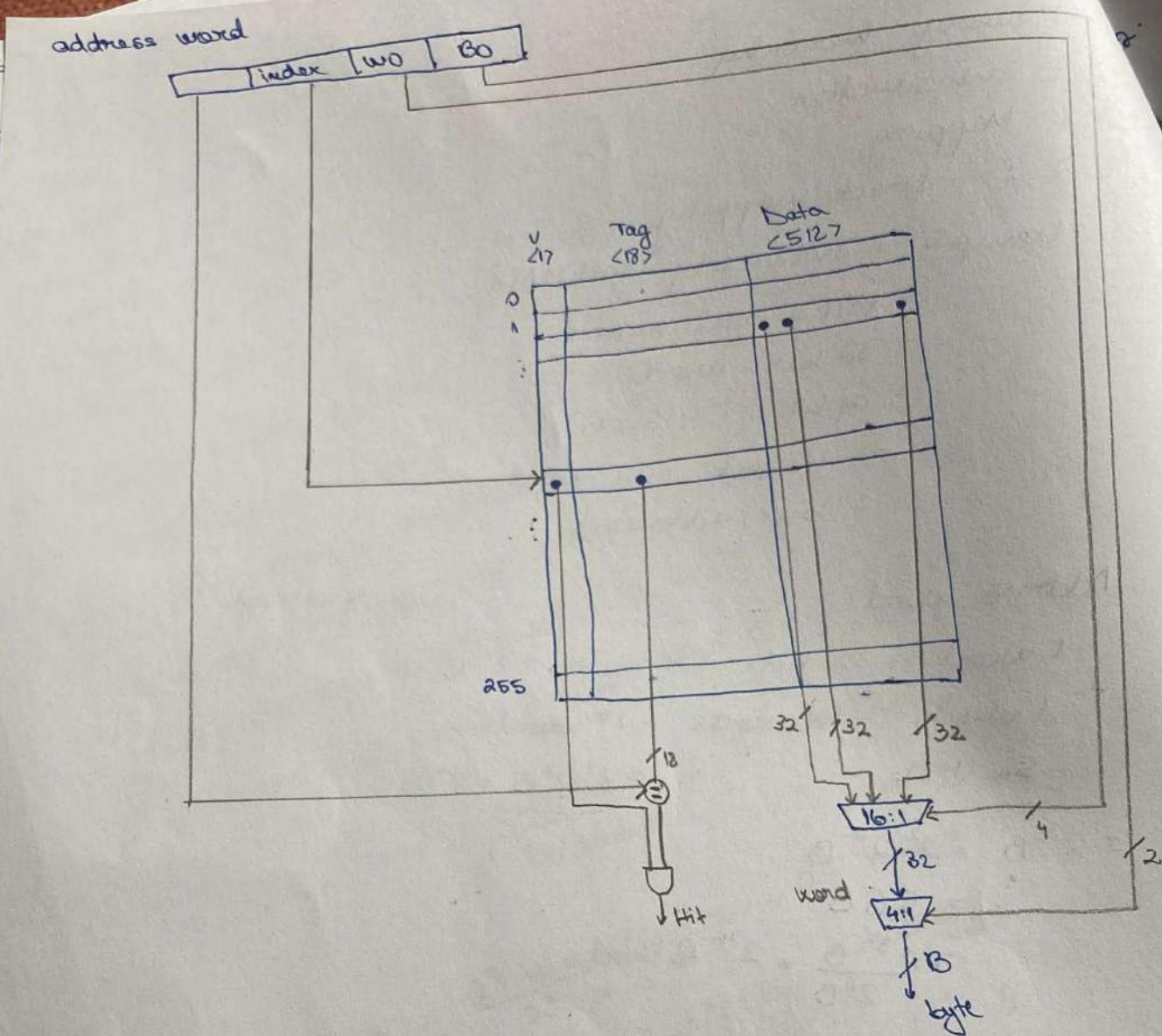
$$\text{Size } (M_1) = 16 \text{ kB} = 2^{14} B$$

$$\frac{2^{14}}{2^{26}} = 2^{-12} \text{ index blocks length}$$

$$\text{Tag length} = M - N = 26 - 8 = 18$$

Address word





3.2.2. Self associative mapping

2^m - blocks at the same cache index

$$M_2 \rightarrow 2^m \text{ blocks} \leftarrow \textcircled{d}$$

$$M_1 \rightarrow 2^m \text{ blocks} = 2^{m-d} \text{ sets} \leftarrow \textcircled{i}$$

$$i = j \bmod 2^{m-1}$$

Exemple :

$$m=5$$

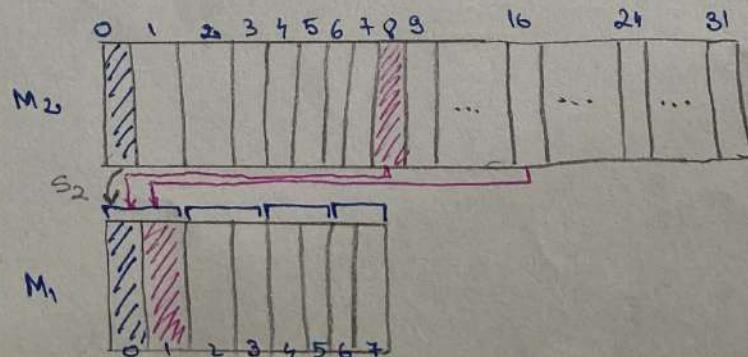
$$m=3$$

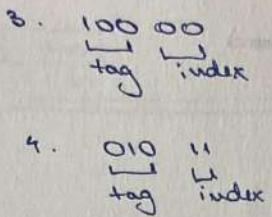
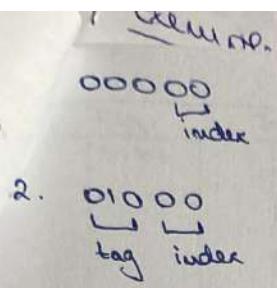
$$d=1$$

$$0 \bmod 2^2 = 0$$

$$8 \bmod 4 = 0$$

$$16 \bmod 4 = 0$$



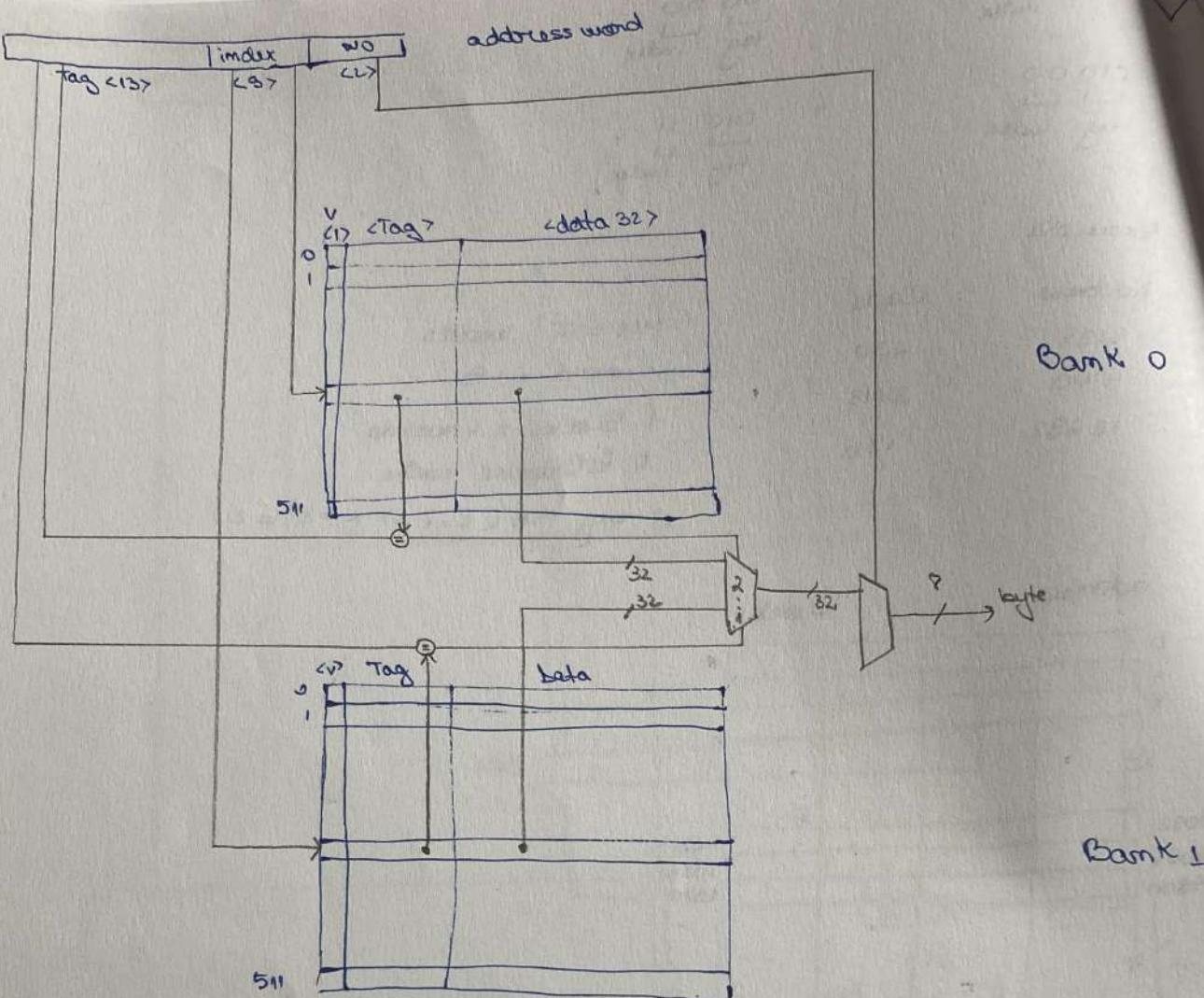


Example:

Address	Code	
3194	230	$MM = 2^{24}$ words
4100	2415	1 word = 1 B
12282	170	1 blocks = 4 words

1 KiBBlocks cache
2-way set ($s=1 \Rightarrow k=2^1 = 2$)

address	blocks
0	0
4	1
8	2
22	
4092	1023
4096	1024
4800	1025
22	
8192	
22	
12288	
12292	
22	
$2^{24}-4$	$2^{22}-1$



$$= 2^8 \text{ B}$$

Bank 1

Tag	Index	w/o
0192	00 00 00 00 00 100	00 00 00 00 00 10
4100	00 00 00 00 00 01 0	00 00 00 00 00 1 00
	00 00 00 00 00 11 0	00 00 00 00 00 1 00

Example:

1 MIPS machine VAX 111780
32 bit word address

cache 2-way set

1 word = 1 byte = 2^0 B

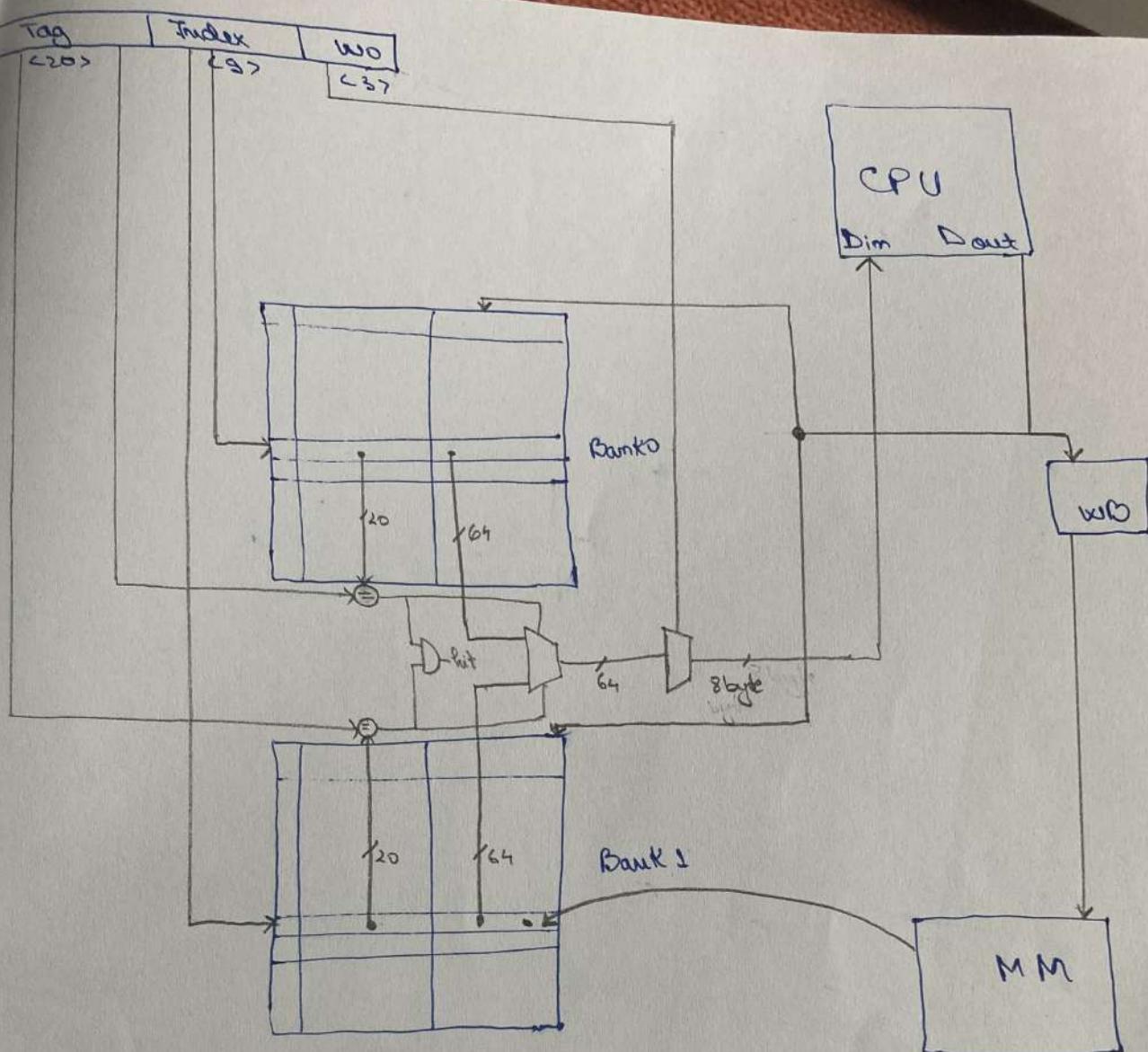
1 block = 8 words = 2^3 words

cache size = 1 KiB

MM size = $2^{32} \text{ B} = 4 \text{ GiB}$

cache size (data) = $2^{10} \text{ blocks} = \frac{2^{10}}{2^3} = 2^8 \text{ sets}$

S10



Example 1

Size (M_2) = 4 GiB - memoria principal

Byte addresses

$$1 \text{ word} = 4 \text{ B} = 2^{20} \text{ B} \quad \text{BO (byte offset)}$$

$$1 \text{ block} = 8 \text{ words} = 2^3 \text{ words} \quad \text{WO (word offset)}$$

4 way SA cache

$$\text{Cache data size } (M_1) = 16 \text{ KiB}$$

a) Address word format = ?

b) Total cache size = ?

c) Cache logic diagram

$$\text{Size } M_2 = 4 \text{ GiB} = 2^3 \cdot 2^{30} \text{ B} = 2^{32} \text{ B}$$

$$\text{Data size } M_1 = 16 \text{ KiB} = 2^4 \cdot 2^{10} \text{ B} = 2^{14} \text{ B}$$

$$1 \text{ block} = 2^3 \text{ words} = 2^3 \cdot 2^2 \text{ B} = 2^5 \text{ bytes/block B} = 2^5 \cdot 2^3 \text{ B} = 2^8 \text{ B}$$

$$M_1 = 2^{14} \text{ B} = \frac{2^{14} \text{ B/cache}}{2^5 \text{ B/block}} = 2^9 \text{ block/cache}$$

$$\text{Bank size} = \frac{2^9 \text{ block/cache}}{2^2 \text{ n.a.}} = 2^7 \text{ block/bank}$$

a)	$\begin{array}{c} <20> \\ \hline \text{Tag} \end{array}$	$\begin{array}{c} <7> \\ \text{Index} \text{WO} \text{BO} \end{array}$	Address word
	31	1211..76543210	

$$b) \underbrace{2^2}_{\text{SA}} \cdot (16 + 76$$

$$\underbrace{2^2}_{\text{SA}} \cdot \underbrace{2^7}_{\text{cache lines/bank}} \cdot (\underbrace{\underbrace{16}_{\text{valid}} + \underbrace{2^4}_{\text{tag}}}_{\text{in plus}} + \underbrace{2^5 \text{ B}}_{\text{data}}) = 30 \cdot 2^9 + 16 \text{ KiB} \approx 18 \text{ KiB}$$

rs 11

Write policies

Cod. coherency (cache coherency)

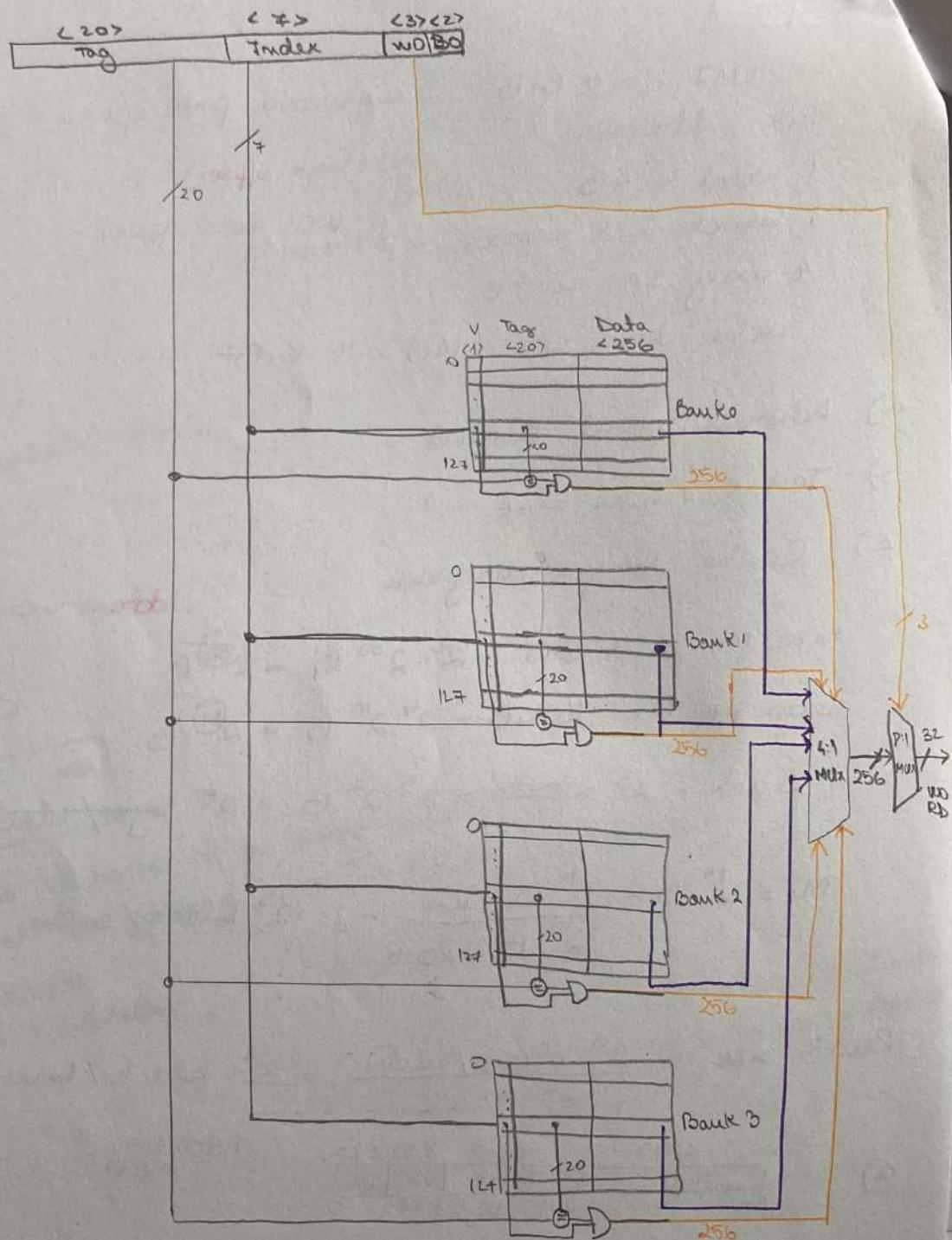
Stat. c)

2) WR

c)

3)

b)



09.05.2023

2.3 Full associative mapping

$K = 2^m$

- ↳ if $s=0 \Rightarrow DM$
- ↳ if $s=M \Rightarrow FA$

Example 2:

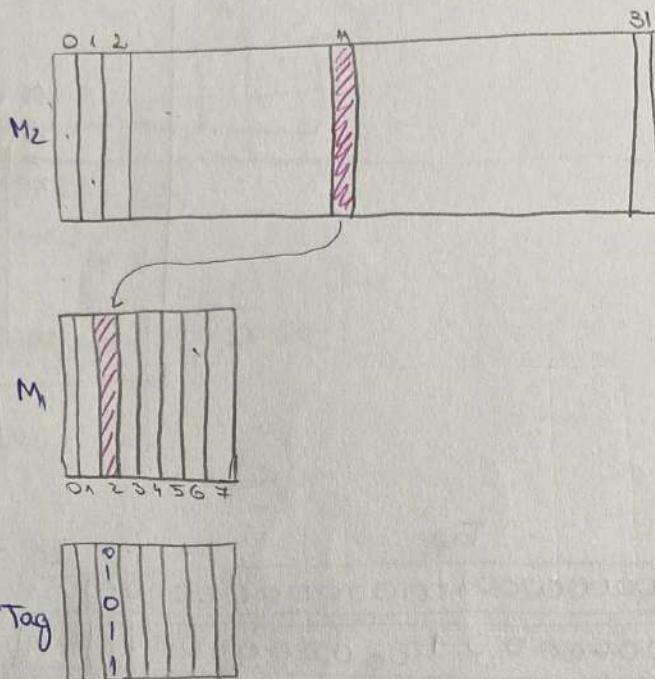
Size (M_2) = 2^n blocks

Size (M_1) = 2^m blocks M_2

$n=5$

$m=3$

" = $\begin{matrix} 0 & 1 & 0 & 1 & 1 \\ \text{Tag} \end{matrix}$



Example 3:

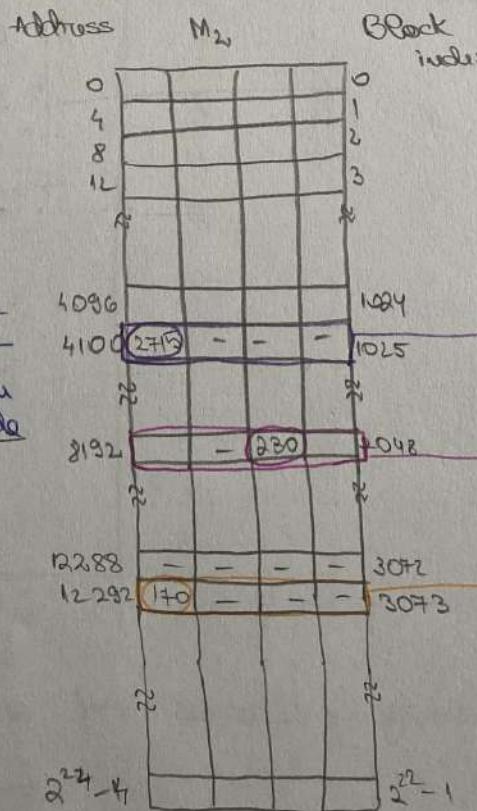
$2^{24} \rightarrow M_2$

1 word $\rightarrow 1B$

1 block = 4 words

M_1 size = 1 kB block

Address	Instruction code
8194	230
4100	2715
12292	170



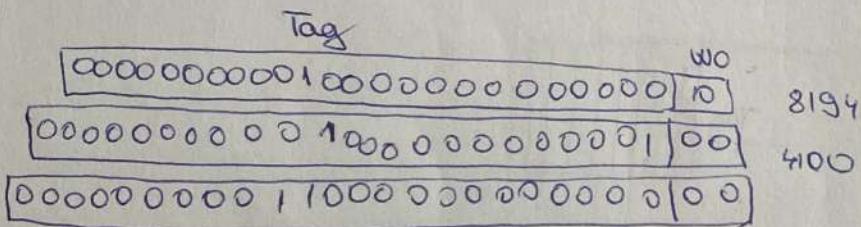
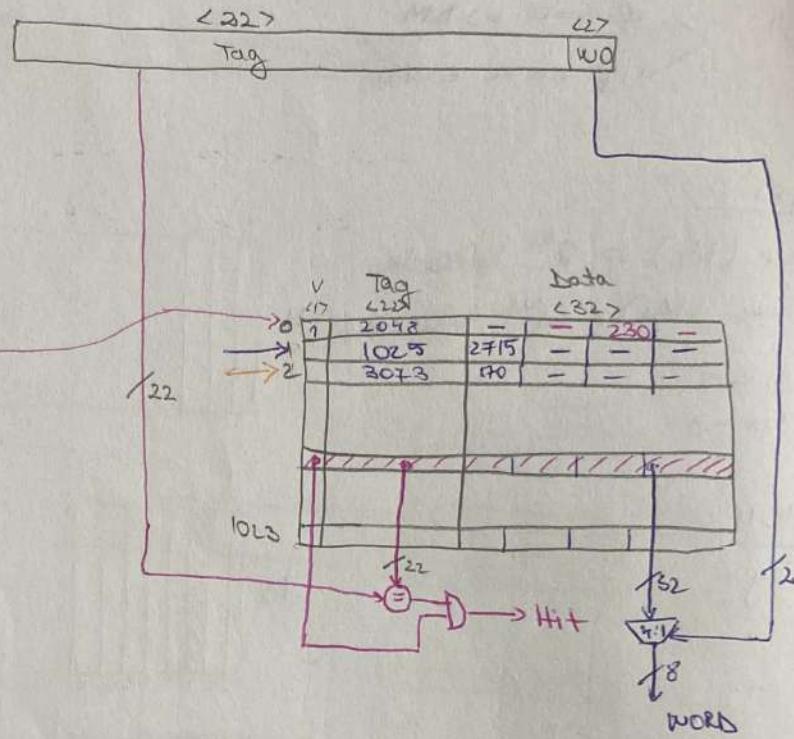
Write policies

Code coherency (cache coherency)

... writes

state
stat
a) write

word
address



09.05.2023

3. Replacement methods

What happens if the cache is full?

SA & FA

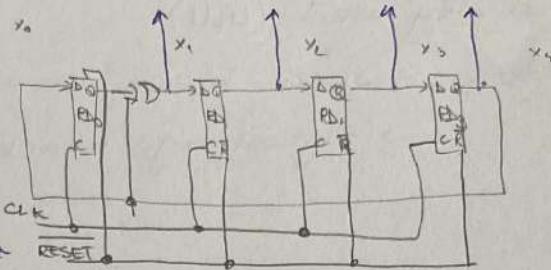
3 Politics

1. Random

Pseudo-random

$$\text{Galois } G(x) = x^4 + x + 1$$

$$\text{Period} = 2^{\deg(G(x))} - 1$$



Linear feedback Shift Register : LFSR

RD0	RD1	RD2	RD3	Dec
1	0	0	0	1
0	1	0	0	2
0	0	1	0	4
0	0	0	1	8
1	1	0	0	3
0	1	1	0	6
0	0	1	1	12
1	1	0	1	11
1	0	1	0	5
				10
				7
				14
				15
				13
				9
				1

Not compliant with the locality principle

May be used for LRU

X
2. FIFO : Not compliant with the locality principle

3. Least recently used (LRU)

Age registers → each block in a set has an age
→ highest age block is replaced

Example 4

4 - ways SR cache

Code	Banks				Age			
	0	1	2	3	0*	1*	2*	3*
M 14	0	0	0	0	0	0	0	0
M 7	14	0	0	0	0	0	0	0
M 6	14	7	0	0	1	0	0	0
M 12	14	7	6	0	2	1	0	0
M 10	14	7	6	12	3	2	1	0
H 6	10	7	6	12	0	3	2	1
H 7	10	7	6	12	1	3	0	2
M 14	10	7	6	12	2	0	1	3
H 7	10	7	6	14	3	1	2	0
M 3	10	7	6	14	3	0	2	1
14	3	7	6	14	0	1	3	2
	3	7	6	14	1	2	3	0

Lec 11

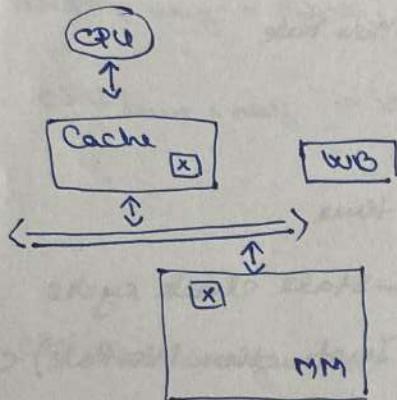
3.4. Write policies

Code coherency (cache coherency)

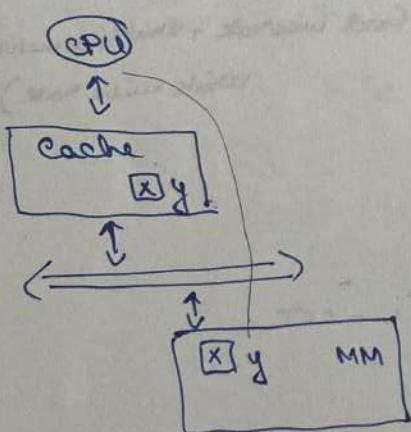
Statistically $\approx 10\%$ are writes

a) Write-through

- write buffer (dimension optimization)



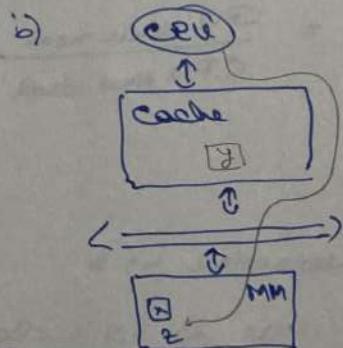
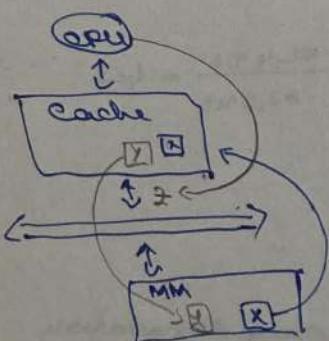
b) Write-back - update upon replacement



What happens in case of a write miss?

[What happens for a read miss? - always allocate]

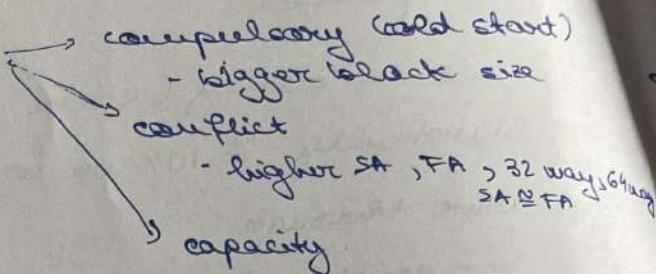
- Write allocate
- Write no allocate


 $b+d - WB + W \text{ allocate}$
 $d+\beta - WT + W \text{ no allocate}$

3.5. Cache memory performance

a) Hit rate \leftrightarrow Miss rate ; $MR + HR = 1$

The 3 main causes of misses



b) Average memory access time (AMAT)

$$AMAT = t_{ac} \cdot Hit\ Rate + (t_{ac} + Miss\ Penalty) Miss\ Rate$$

$$= t_{ac} + Miss\ Rate \cdot Miss\ Penalty$$

t_{ac} = cache access time (1cc - 2cc)

c) CPU time = Clock cycles · Clock cycle time

$$\text{Clock cycles} = \text{Execution clock cycles} + \text{Stall clock cycles}$$

$$= IC \cdot CPI + IC \cdot (\text{Memory access per instruction} \cdot \text{Miss Rate}) \cdot CCT$$

$$\text{CPU time} = IC \cdot CPI_{\text{ideal}} + (\text{Memory access per instruction} \cdot \text{Miss Rate} \cdot \text{Miss Penalty}) \cdot CCT$$

$$\text{Stall clock cycles} = (\frac{\text{Reads per instruction}}{\text{Reads per instruction}} \cdot \text{Read miss rate} + \text{Writes per instruction} \cdot \text{Write miss rate}) \cdot \text{Miss penalty}$$

Example 1:

$$CPI_{\text{ideal}} : 8,5 \text{ cc} \quad \text{Miss penalty} : 6 \text{ cc}$$

$$\text{Memory access per instruction} : 3 \quad t_{ac} = 5 \text{ ns} = CCT$$

a) AMAT = ?

$$\text{Miss penalty} = 6 \text{ cc} = 6 \cdot 5 \text{ ns} = 30 \text{ ns}$$

$$AMAT = 5 \text{ ns} + 0,11 \cdot 30 \text{ ns} = 8,3 \text{ ns}$$

b) CPU time ideal = $IC \cdot 8,5 \cdot 5 \text{ ns} = IC \cdot 42,5 \text{ ns}$

$$\text{CPU time real} = IC \cdot \underbrace{(8,5 + 0,11 \cdot 6)}_{\text{cc}} \cdot 5 \text{ ns} = IC \cdot 52,4 \text{ ns}$$

$$\frac{\text{Performance ideal}}{\text{Performance real}} = \frac{\text{CPU time real}}{\text{CPU time ideal}} = \frac{IC \cdot 52,4 \text{ ns}}{IC \cdot 42,5 \text{ ns}} = 1,23$$

Example 2:

Cache : 64 KiB

CCT = 20 ns

Memory access / interaction = 1,3

CPI ideal = 1,5 cc

SA \Rightarrow additional MUX \Rightarrow 8,5% clock frequency degradation

SA

(1) Which is better?

1 way SA or 2 way SA

1 way SA

2 way SA

MR 3,9%

3%

MR 200 ns

200 ns

$$a) \text{ATMAT}_{\text{1 way}} = 20 \text{ ns} + 0,039 \cdot 200 \text{ ns} = 27,8 \text{ ns}$$

$$\text{ATMAT}_{\text{1 way}} = 20 \text{ ns} \cdot 1,085 + 0,03 \cdot 200 = 27,7 \text{ ns} \quad \checkmark$$

$$b) \text{CPU time 1way} = \text{ic} \cdot (1,5 + 1,3 \cdot 0,039 \cdot \text{Miss Penalty}^c) \text{ CCT}$$

$$= \text{ic} \cdot (1,5 \cdot 20 \text{ ns} + 1,3 \cdot 0,039 \cdot \underbrace{\text{Miss Penalty}^c}_{20 \text{ ns}})$$

$$= \text{ic} (30 \text{ ns} + 10,14 \text{ ns}) = \text{ic} \cdot 40,14 \text{ ns} \quad \checkmark$$

$$\text{CPU time 2way} = \text{ic} (1,5 + 1,3 \cdot 0,03 \cdot \text{Miss Penalty}^c) \text{ CCT}$$

$$= \text{ic} [(1,5 \cdot 20 \cdot 1,085) \text{ ns} + 1,3 \cdot 0,03 \cdot 200 \text{ ns}]$$

$$= (32,55 \text{ ns} + 7,8) \text{ ic}$$

$$= \text{ic} 40,35 \text{ ns}$$

Write back

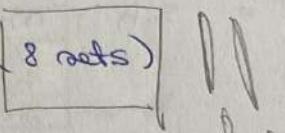
- Update upon replacement
- 1 dirty bit

Miss - allocation

Exemple de problème (Examens)

1.

Cache (instruction) 4-way SA (8 sets)



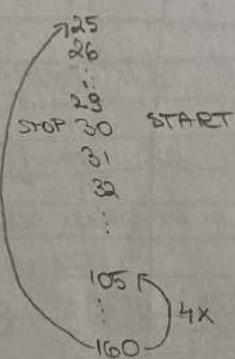
- 1 block = 4 words
- 2^{16} words $\rightarrow M_2$
- Addressable unit \rightarrow word
- Cache : initially empty
- $t_{ac} = 30 \text{ ns} = \text{clock cycle time}$
- Miss Penalty $= 480 \text{ ns}$ (c.r.)
- LRU replacement

a) Memory address format =?

Total cache capacity, 1 word = 32 B

b) Hit ratio and AMAT,

- execută un program dat ca adresă a securuților :



29 = 000111|011

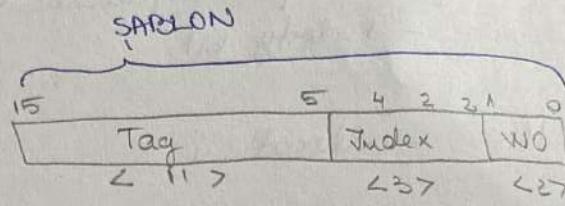
Tag = N_H mod nr. de seturi

=

Ques 13

Relevance:

- a) we can byte offset (BO)
 due 1 block = 4 words = $2^2 \rightarrow$ WO
 $2^{16} \rightarrow$ Address word
 2^{16} words $\rightarrow M_2$



$$S.A : 4 = 2^2$$

$$\text{Total Cache} = S.A \cdot \underbrace{\text{no. setwi / bank}}_8 \cdot \left[\underbrace{(1 + 1)}_{\text{valid}} + 2^2 B \right]$$

Blocks/cache

$$1 \text{ word} = 32 \text{ bits} = 4B = 2^2 B$$

instruction

$$\begin{aligned} \text{Total Cache} &= 4 \cdot 8 \cdot (12 + 4) \\ &= 4 \cdot 8 \cdot 16 \\ &= 2^2 \cdot 2^3 \cdot 2^4 \\ &= 192 B \end{aligned}$$

b)

M_2

Address	Block number
0 - 3	0
4 - 7	1
8 - 11	2
12 - 15	3
16 - 19	4
20 - 23	5
24 - 27	6
28 - 31	(7)
32 - 35	8
36 - 39	9
40 - 43	10
44 - 47	11
48 - 51	12
52 - 55	13
56 - 59	16
60 - 63	15
64 - 67	16
68 - 71	14

72 - 75	18
76 - 79	19
80 - 83	20
84 - 87	21
88 - 91	22
92 - 95	23
96 - 99	24
100 - 103	25
104 - 107	26
108 - 111	27
112 - 115	28
116 - 119	29
120 - 123	30
124 - 127	31
128 - 131	32
132 - 135	33
136 - 139	34
140 - 143	35
144 - 147	36
148 - 151	37
152 - 155	38
156 - 159	39

a) WWT
b) WB

		$10^8 \times 0.1$
09		$160 - 163$
15		$32 - 35$
1		$36 - 39$
1		$40 - 43$
1		$44 - 47$
1		$48 - 51$
1		$52 - 55$
0	10	$56 - 59$
1		$24 - 27$
1	0	$28 - 31$
4		$80 - 159$

23.05.2022

V	Tag	
0	1	2
1	1	2
2	0	2
3	0	2
4	1	2
5	1	2
6	0	2
7	1	1

B1

$160 - 163$	40
-------------	----

2846

V	Tag	
0	1	3
1	1	3
2	1	3
3	1	3
4	1	3
5	1	3
6	1	3
7	1	2

B2

V	Tag	
0	0	4
1	0	4
2	0	4
3	0	4
4	0	4
5	0	4
6	0	4
7	1	3

B3

$$(IM) \times 31$$

N

Cod

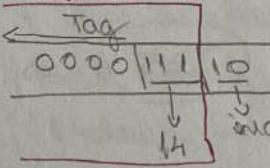
fig. num. de la

ic

Prima instr. la adresa 30 $\Rightarrow 1M$ (1 miss)

Alocă blocul, blocul se află în tabel 28-31

Tag: cîntul împărțitii lui Block number la



\rightarrow ŞABLONU DE LA ÎNCEPUT

8 module 8

8 \Rightarrow Tag = 1

9 module 8 - 1

Ajung la 23 pe ! 23 modul 8 = 7 \Rightarrow măduse la indexul 7

Ajung la 156 în blocul 39 \Rightarrow 39 modul 8 = 7

Dacă cîntul e plin, la indexul 7 dar LRU deci replosez indexul 7

dim BO

$$\frac{156}{M} = \frac{157, 158, 159}{4} \Rightarrow (1M) \times 1$$

La $160 \cdot (1M)$ în blocul 40, 40 modul 8 = 5 deci rep. indexul

dim BO

02.05.2023

Avei executat tot de la 30 la 160

\Rightarrow sare însoțită la 105 care e în cash

105 - 160 cash-ul are instrucțiunile

$$\Rightarrow \begin{pmatrix} 56 \\ 4 \end{pmatrix}^+ = \begin{pmatrix} 224 & 4 \\ 0 & M \end{pmatrix}$$

O fac de 4 ori (bucle)

\rightarrow astăzi valoari sunt deja în tabel de la 160 la 105

Deci avei executat și bucle multe

sare însoțită la 25

25 lipsește \Rightarrow 1M

replasez pe călă mai vechi

~~26, 27, 28~~ - sunt \rightarrow 2H

~~26, 27~~

28 lipsește 1M

29, 30 \rightarrow 2H

$$\Rightarrow \begin{pmatrix} 1M \\ 1H \end{pmatrix} \begin{pmatrix} 1M \\ 3H \end{pmatrix} * 32 \begin{pmatrix} 1M \\ 224H \end{pmatrix} \begin{pmatrix} 1M \\ 2H \end{pmatrix} \begin{pmatrix} 1M \\ 2H \end{pmatrix}$$

deci aveam

$$\begin{array}{r} 131 + \\ 224 \\ \hline 231 + \\ 131 \\ \hline 361 \text{ acese} \end{array}$$

36M

$$\begin{array}{r} 101 + \\ 224 \\ \hline 325 \end{array}$$

325H \Rightarrow 361

$$\text{Hit ratio} = \frac{325}{361} = 0,9002 \approx 90\%$$

continuare: calcul AMAT

$$AMAT = 30 \text{ ns} + 0.9 \cdot 480 \text{ ns}$$

Exemplu 2.

Sist. de calcul load store cu un cache
cu un cache
cu un cache

atunci când rulează ceea:

10% Data Miss Rate

8% Instruction Miss Rate

20% Load Store instructions

15% of the blocks are dirty at any time

CPI = 2 cc dacă cash-ul este perfect

CPI ideal = 2 cc

I C = 10 000

Clock frequency = 4,5 GHz

Miss Penalty este determinată de faptul că un acces la bus revine la 10 cc

1 Bus access = 10 cc

Bus width = 4 words

1 block = 16 words

Politica de tip Write Back + Write Allocate (entire block)

a) AMAT = ?

b) CPU time = ?

$$AMAT = t_{ac} + \text{Miss Rate} \cdot \text{Miss Penalty}^+$$

$$\text{Clock cycle time} = \frac{1}{4,5 \cdot 10^9} = \frac{1}{4,5 \cdot 10^9} \approx = \frac{1 \cdot 10^{-9}}{4,5} = 0,22 \text{ ns}$$

Clock freq.

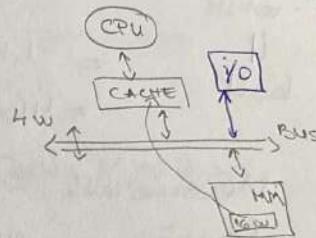
$$\text{Miss Penalty}^+ = \text{Clock cycle time} \cdot \text{Miss Penalty}^c$$

Miss Penalty^c

$$16 \text{ words allocation} \Rightarrow \frac{16}{4} = 4 \text{ Bus accesses}$$

Bus width

$$10 \text{ cc} \Rightarrow 4 \cdot 10 \text{ cc pt. } \boxed{\text{alocare}}$$



Cours 13

Read Miss Penalty = allocation + update

40 cc ← allocation penalty
update penalty

Read Miss Penalty = Allocation Penalty + Percentile de dirty, Update Penalty
0,15

writen

$$\begin{aligned} \text{||} &= 1,15 \cdot 40 \\ \text{back} &= 46 \text{ cc} \end{aligned}$$

Mecanismul este asemănător cu Writing Miss Penalty
(doar nu există Back nu sunt egale)

$$\text{Write Miss Penalty} = 1,15 \cdot 40 = 46 \text{ cc}$$

$$\Rightarrow \text{Read Miss Penalty}^c = 46 \text{ ms}$$

$$\Rightarrow \text{Miss Penalty}^t = 46 \cdot 0,22 = 10,21 \text{ ns}$$

$$\text{AMAT} = 0,22 \text{ ns} + \dots$$

10% Data Miss Rate - numai din când în când (la read store)

8% Instruction Miss Rate - de fiecare dată

$$\text{Miss Rate} = 0,08 + 0,2 \cdot 0,1 = 0,1$$

$$\downarrow \quad \downarrow$$

$$L/S \quad \text{Instruction} \quad \text{Instruc. Miss Rate}$$

$$\text{AMAT} = 0,22 \text{ ns} + 0,08 \cdot 10,12 \text{ ns} = 1,23 \text{ ns} = 0,22 + 0,0833 \cdot 10,12 = 1,06$$

$$\text{Miss Rate} = \frac{\text{Misses Per Instruction}}{\text{Mem. access per instruction}} = \frac{0,1}{1 + 0,2}$$

$$\downarrow \quad \downarrow$$

$$pt. citire \quad pt. acces la memorie$$

$$instructione \quad pt. date$$

- 1005/13
 a) WWT
 b) WWB

23.05.2023

$$\text{CPU time} = IC \left(\text{CPI ideal} + \text{Memory accessed per instr.} \cdot \frac{\text{Miss Rate}}{\text{Miss Rate Penalty}} \right)$$

• Clock cycles time

$$= 10000 (2 + \text{Mem. acc. per instr. Miss Rate} \cdot 46) \cdot 0,22 \text{ ns}$$

lifey - Up to memory access per instr. • Miss rate => we take access to memory.

penalty

count Miss

$$\Rightarrow \text{Mem. acc. per instr. Miss Rate} = \text{Misses per instr.}$$

=

$$= 1 \times \text{Instruct. Miss Rate} + \text{L/S Instruction. Data Miss Rate}$$

$$= 1 \cdot \text{IMR} + 0,12 \cdot \text{DMR} = 0,1$$

$$\begin{aligned}\text{CPU time} &= 10000 (2 + 4,6) \cdot 0,22 \\ &= 6,6 \cdot 0,22 \cdot 10000 \\ &= 66 \cdot 22 \cdot 10 \\ &= 14520 \text{ ns} \\ &= 14,520 \mu\text{s}\end{aligned}$$

Example 3

a) % BUS Used

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

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

$$\text{Hit rate} = 90\%$$

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

In case of miss a block is allocated

Block allocation

$$\text{CPU reference frequency} = 10^2 \text{ words/sec}$$

- 30% writes

35% of blocks are dirty

Write allocate for a write miss

- a) cache-ul este write through
 b) cache-ul este write back

$$10^8 \text{ words/sec} \cdot 0,1 \cdot \frac{\text{Miss Penalty}^b}{\text{Miss Rate}}$$

$$\text{Miss Penalty}^b = 0,7 \text{ Read miss penalty}^b + 0,3 \text{ Write miss penalty}^b$$

$$\text{Read miss penalty}^b = 0,35 \cdot \underbrace{\frac{4 \text{ words}}{2 \text{ words}}}_{\text{update writes}} + 1 \cdot \underbrace{\frac{4 \text{ words}}{2 \text{ words}}}_{\text{reads}} = 1,35 \cdot 2 = 2,7$$

$$= \text{write miss penalty}$$

23.05.2023

Cours 13

a) WT

b) WB

$$\begin{aligned}
 & 10^8 \times 0,1 \times (\% \text{ Reads} \times \text{Read Miss Penalty}_{WT} + \% \text{ Writes} \times \text{Write Miss Penalty}_{WB}) \\
 & \text{Read Miss Penalty}_{WB} = \text{Write Miss Penalty}_{WB} \\
 & = 0,35 \times \frac{4 \text{ words}}{2 \text{ words}} + \frac{4 \text{ words}}{\underbrace{2 \text{ words}}_{\text{Bus Reads}}} \\
 & \quad \underbrace{\text{BUS Writes}}_{\text{BUS Accesses}} \\
 & = 1,35 \times 2 \\
 & = 2,7
 \end{aligned}$$

$$10^8 \times 0,1 \times 2,7 = 2,7 \times 10^7$$

$$= 2,7 \%$$

2,3

$$a) 10^8 \times 0,1 \times (\underbrace{\% \text{ Reads} \times \text{Read Miss Penalty}_{WT}}_{2,3} + \% \text{ Writes} \times \text{Write Miss Penalty}_{WT})$$

$$+ 10^8 \cdot 0,3 \times \underbrace{0,9}_{\text{Hit}} \times \text{Write Hit Penalty}$$

$$\text{Read Miss Penalty}_{WT} = \frac{4 \text{ words}}{\underbrace{2 \text{ words}}_{\text{Allocate}}} = 2 \text{ Bus reads}$$

$$\text{Write Miss Penalty}_{WT} = \frac{4 \text{ words}}{\underbrace{2 \text{ words}}_{\text{Allocate}}} + 1 \text{ Bus writes} = 3 \text{ Bus Acc.}$$

$$\text{Write Hit Penalty}_{WT} = 1 \text{ Bus write} = 1 \text{ Bus Access}$$

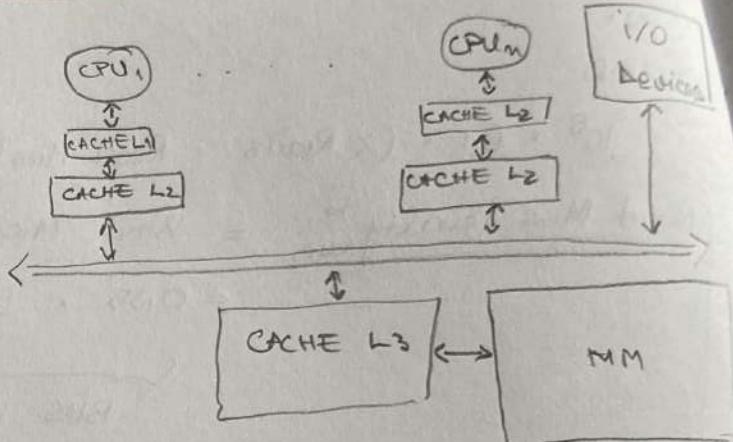
$$0,7 \cdot 2 + 0,3 \cdot 3 = 2,3$$

$$\Rightarrow 10^7 \cdot 2,3 + 10^7 \cdot 2,7 = 5 \times 10^7$$

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

3.6 Reducing miss penalties and increasing cache performance

- victim caches
- multiple cache levels



Example 1

$$CPI_{ideal} = 1 \text{ cc}$$

$$\text{Clock rate} = 4 \text{ GHz} \Rightarrow \text{Clock cycle time} = \frac{1}{4 \cdot 10^9} \text{ s} = 0,25 \text{ ns}$$

MM access time = 100 ns (including all misses handling)

Miss rate per instruction = 2% (at the primary cache)

L₂ cache = 5 ms access time

New miss rate = 0,5%

Miss Penalty MM = 400 ns

$$CPU_{time\ 1-level} = IC \times (1 + 0,02 \times 400) = IC \times 9cc \times 0,25 \text{ ns}$$

Mem. access per instruction * Miss Rate = Misses per instruction

$$CPU_{time\ 2-level} = IC \times (1 + 0,02 \times 20 + 0,05 \times 400) \times 0,25 = 3,4$$

$$\frac{\text{Performance 2-level}}{\text{Performance 1-level}} = \frac{CPU_{time\ 1-level}}{CPU_{time\ 2-level}} = \frac{9cc}{3,4} = \underline{\underline{2,65}}$$

Speedup w/
odus de el doble
lvl.

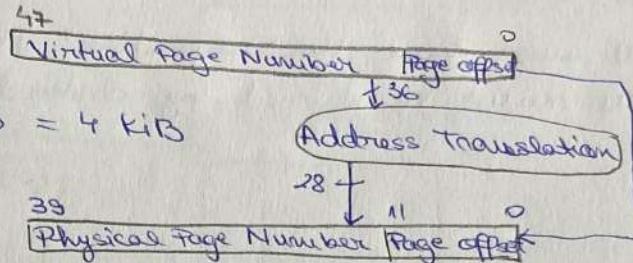
ARM V8 Address word \rightarrow 64 bit

MS 16 bits are unused

Virtual Address

$$\text{page size} = 2^{12} \text{ B} = 4 \text{ KiB}$$

Physical address



$$2^{48} \text{ B} = 256 \text{ TiB} \quad \text{Virtual Mem Space size}$$

$$2^{40} = 1 \text{ TiB} \quad \text{Physical Mem size}$$

Design choices \rightarrow huge miss penalty

① Large pages : 1 KiB - 64 KiB

Embedded : 1 KiB - 4 KiB

Desktop, Servers : 32 KiB - 64 KiB

② Full associative mapping

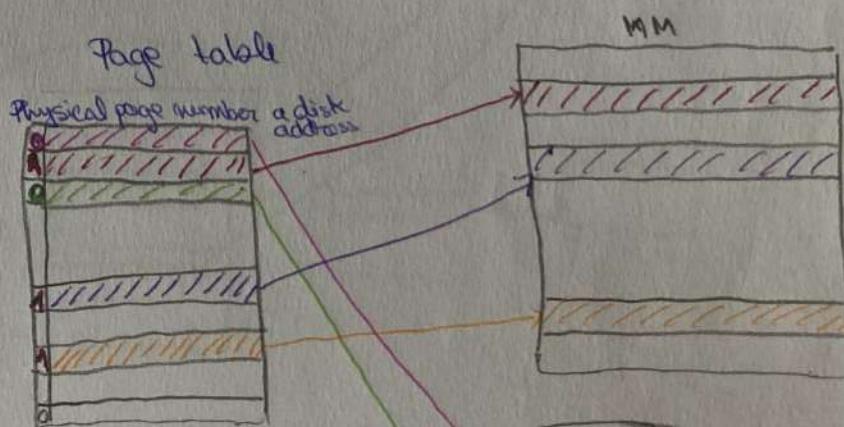
- implemented through page tables

WB is the only choice

Page table register

Virtual address

VPN PO



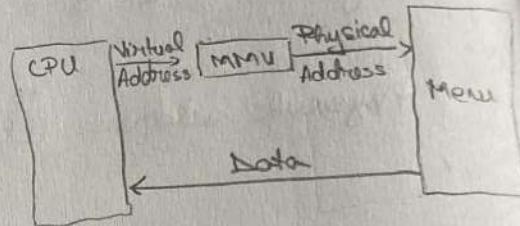
Page table + PO + Internal register = Page state

Secondary stage

4. Virtual memory

4.1. Motivation

1. CPU must have the "impression" of ∞ size memory
2. CPU must be shared by distinct processes/virtual mach.



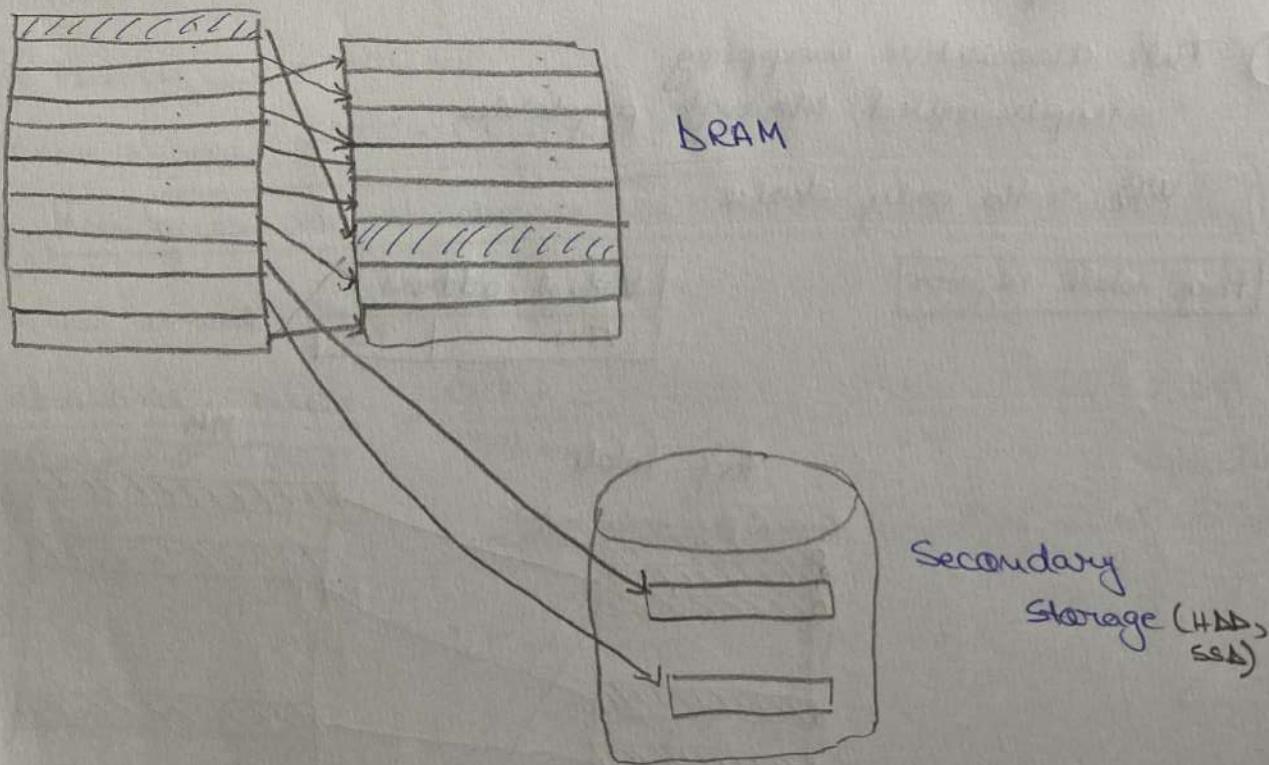
Page table
Translation Table
Base Register (TR)

Locality principle

Blocks $\xrightarrow{\text{Segments : dimensione variabili}}$
 $\xrightarrow{\text{Pages : dimensione fissa}}$

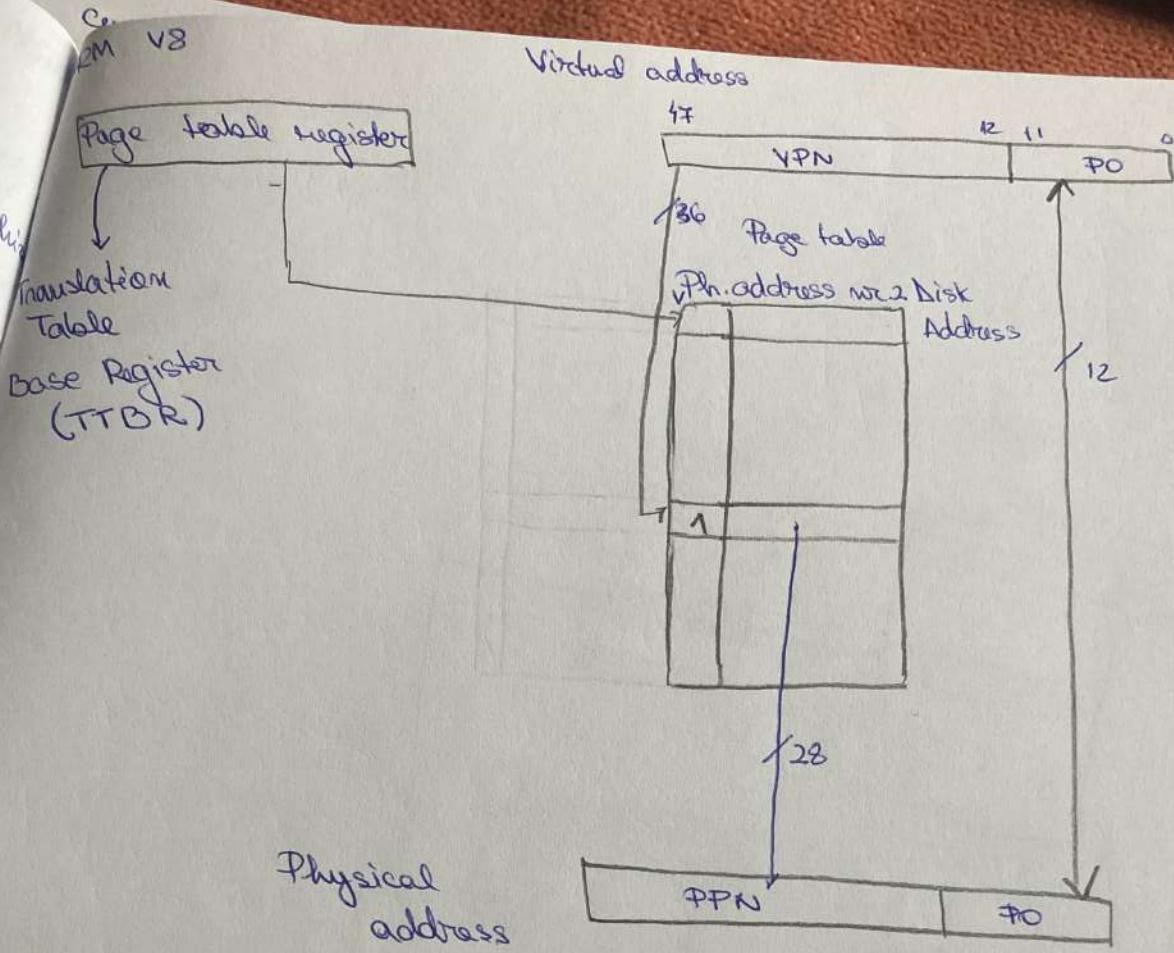
Virtual Memory

Physical Memory



32 bits
- 32 bit

>

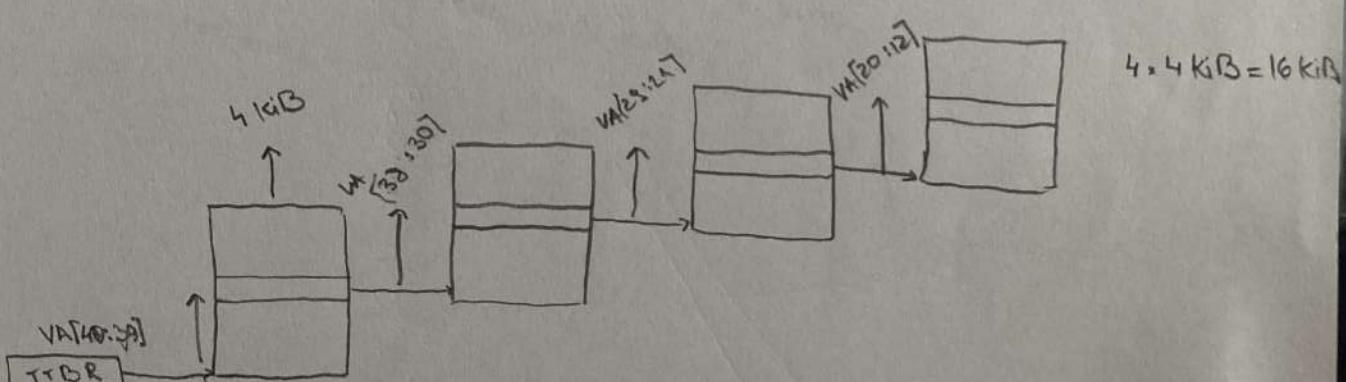


4.2 VMM performance optimization

ARM V8 page table size = $2^{36} \times 2^2 B = 2^{38} B = 0,25 TiB$
 \Rightarrow Page table size reduction

- paged tables
- hierarchical page tables

$$2^2 \cdot 2^{10}$$

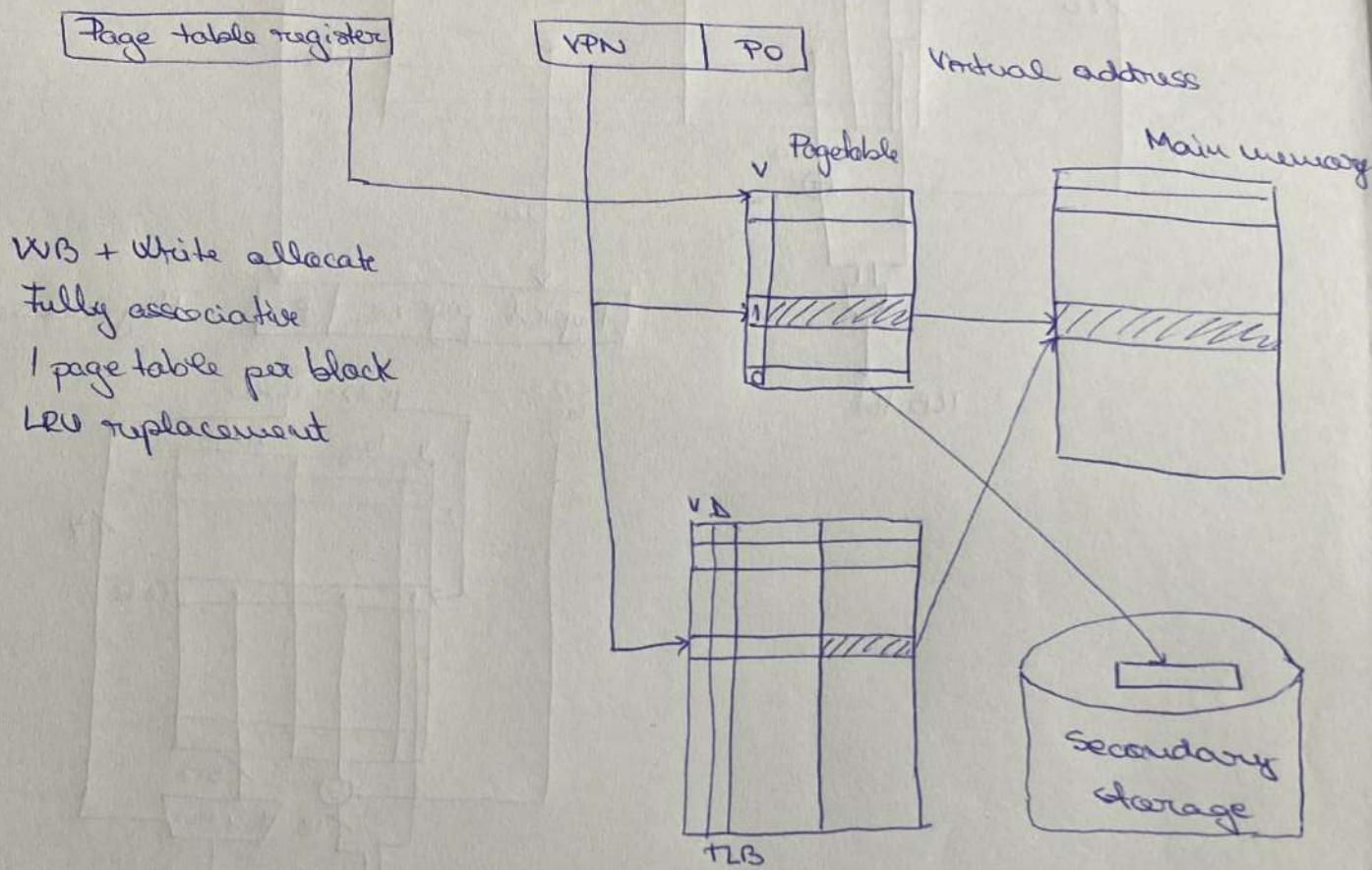


- ② Access time \rightarrow Cache for addresses

Cours 14

4.2 VM optimization

Caching the page table \rightarrow Translation Lookaside Buffer (TLB)



Example 1:

Intrinsity's fast Math

TLB 16-512 entries

block 1-2 page table entries

hit time > 1 ns

miss penalty = 10_n-100_D ns

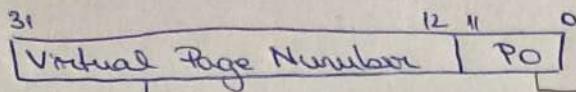
Miss rate 0,1 - 1%

FA

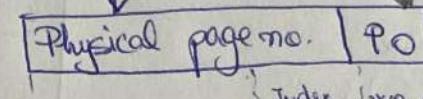
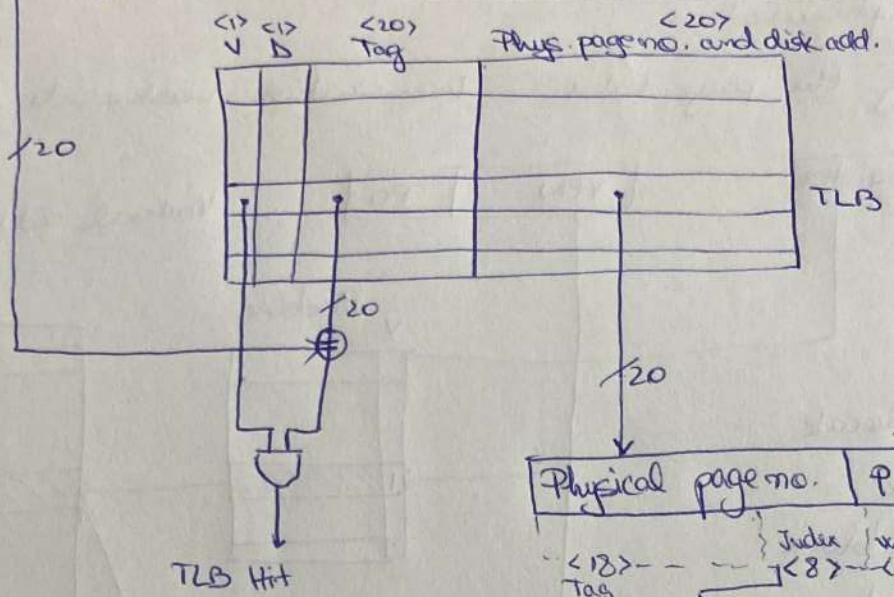
Virtual address - 32 bits

Physical address - 32 bits

Page size = 4 kB



Virtual address



Physical add.

